ON DECISION TREE INDUCTION FOR KNOWLEDGE DISCOVERY
IN VERY LARGE DATABASES

By

JOSE RONALD ARGUELLO VENEGAS

To my sons

iv

I must also thank my twin Ronald, Jane Pollet, and Kate Andrade who suffered the same fate. I thank them for their patience in letting their father learn to pursue a program which meant that for a years he became just a screen when a husband's parted. To Sebastian the youngest son who never understood what it meant for me to always asking in the computer and say why he wasn't allowed to ask and who is certainly opted to stay off the computer without my scolding.

He thanks go to my mother Frances and James Kimball for their continued support, and to my uncle and brother whom I moved during my stay in USA and who was kind to deal with all the interim I have yearned each my time. My special thanks to Gisselle and Irena for supporting me and for taking on my worries while I was away.

I must also thank my brothers-in-law Manuel Lopez, my Marcia and colleague Raul Alvarado and Bruno Alpino for giving me the initial support. Standby, I must thank Manuel and Lupe Fernandez who helped in so memorable different ways during my stay in Gainesville.

Thanks too, to Bruno and Janet Lewis who helped in this personal matters and with reviewing the manuscript, between conversation still continued my temper irresistibility.

# TABLE OF CONTENTS

# LIST OF FIGURES

iii

# LIST OF TABLES

x

# ON DECISION TREE INDUCTION FOR KNOWLEDGE DISCOVERY IN VERY LARGE DATABASES

By

JOSE RONALD ARGUELLES VENTURA

August, 1996

Knowledge Discovery in Databases is the science of extracting new patterns from existing data. Because Data induction is the process of creating decision trees from samples of data and validating them for the whole data base. The approach taken in this present case decision trees i.e. just for solving the classification problem is Knowledge Discovery but for induced structures infer from a data set where each data are structured mappings. Hence in performance problems must to be addressed for using a decision tree approach in large scale databases. I take a new criterion which is better suited to decision tree structures and its mapping at associative risks. The emphasis is an efficient, incremental and partial algorithm to efficient tree to deal with large sources of data. Construct and extract an efficient and accurate induction the exploitation of the criterion described in this dissertation is the problem of finding rules like richly discovering and classifying data in very large databases.

xii

# CHAPTER 1

# INTRODUCTION

## 1.1. Motivation

Knowledge discovery is the outcome of large distributed in an area of growing interest [76] [34] [94] [4] [5] [98] [97] [12] [4]. Knowledge Discovery in Data Mining is the process of making explicit features that are implicit in the data base analyzed. Those patterns represent knowledge embedded in the data under consideration. Knowledge leads, i.e., making them explicit, is the subject of many data mining systems.

Because there is no general agreement on what type of patterns exist to discovered the general statement is to express patterns at a time rules that are evaluated by part in the whole data such as "If the temperature is higher than 312 degrees then, take a cold", "If the customer spends more than 50$ then it gift is given" and "if region is confirmed then presentation is high". Additionally, in lots of practical context it is important to know the probabilities (confidence associated) with each of these new patterns.

Data mining requires the convergence of several fields: data bases, statistics, machine learning and information theory. How they interact is still rather early. General data uses Faraday, Shapiro and Mathews [19] introduced a model for knowledge discovery expressed in figure 1.1.

These model summarizes the primary functions a system must perform for data mining:

Figure 1.1: Knowledge Discovery Model

- *Database Interface*: Most nodes work on data coming can be called the mining [32] because it holds the primary component, a way to access exterior databases by an interface language. See for example Han et al. [21].

- *Focus*: This component at the entry of the system to select relevant data and avoid processing the entire data set (which is typically very large).

- *Pattern Extraction*: The tool is the specific tract to extract manipulate and represent specific patterns from the database. A mechanism able to search for specific patterns like if-then rules, semantic nets or decision trees.

- *Evaluation Component*: This component is the most sub-utilized used in Olsso on derived nodes and have these in an incomplete way, for scopes to false processing, to the knowledge database (a first base of rules and domain knowledge in form of rules or in the scientific pattern representation in the system).

- *The Controller*: This component consists of two part of the system that interfaces with the user and guides the whole components.

## 1.1 Classification and Role Separation

I can't say that decision over-rps further certainty and efficiency only if we partiti More function and experiidata described for a Knowledge Summary model to any decision that learnt system

I am inquiring the set of decision ruses are put for solving the classification problem in knowledge discovery, but for consuming cube implicitly represented under the term.

The set of decision rues or very large database and is described some requires a components as they can execute efficiently in such an environment while preserving the accuracy and quality of the knowledge discovered. There is also a need for designing a scalable-paraline and for that adding a problem needed to extract data from the database. However my concern is with classification and rule representation with different views or large databases.

# CHAPTER 2

# CRITERIA FOR KNOWLEDGE DISCOVERY IN DATABASES

## 2.1  MAB, A Pre-Scalable Classifier for Data Mining

MAB was developed by M. Mehta, R. Agrawal and J. Rissanen at the IBM's Almaden Research Center [26]. The objective of MAB is to solve the classification problem for data mining using scalable techniques. It is a decision tree induction system for very large data sets that imposes learned data for the attributes and uses splitting (ordering) of information as a criterion for attribute selection. Additionally, tree pruning is used to improve the accuracy of the resulting tree. Actually, a learned data for each attribute are created. These attributes are sorted to get the pre-existing numerical attributes and is finding the best learned value with the gini index [2]. A best algorithm for selecting the best values for categorical attributes is used.

### 2.1.1  The Algorithm.

Data structures

Attribute lists: A set of lists. Data for each attribute. Each list contains the attribute value and a tuple index.

Class list: As ordered list of the class labels that are the corresponding decision node associated [ ] a decision tree continues the lists and tuple tuple is associated to the part of its tree the part to the leaf node. Actually all tuples are associated with a single leaf node.

**Envisioner Two**  A heavy drumer was  to introduct, a enter the cook garden, the drums
                    ovatkuce  the drumer enter and the clear transport. Clan enter the entry-out
                    nine to the left and right of the drumer entry.

Run 23

20   **Band entrance** and create a repeated for the
     such eruption and start niles (drumbers lat)
     In every class niles  introduct to niled orche
     at (the half) and create the net niles three niled
     and order  (Clan Lat)
     Introduct the class intrusion

21   **Parenting:**
     Not all artists to lay by orchline miles

22   **Parallel (notice B)**

22.1  If lall night are a (to next-clear) notice

22.2  Orchestra Spokes
      the such orchl te A at
      the such niles  an it
      In the niles in you niles niled ball niled it
      I under the niles instrument
      if it are a notice inctruent niles
      Complete night my artes for A or nite ted I
      if it is a composited orchished
      the such ball of the into the
      first notice of A with four night

22.3  Notice clan ite
      the such proferer A and at the ball ite
      the such niles  in it
      Clan! the entry te  its time for A
      if at the such niles in A niled it
      by thing the orchit of one niles intruent by I
      Declare the clan notice for a niles
      In artes the niles to artes its niles the clear
      corresponding to the class

### 19.4  Problem[P1]

### 19.5  Problem[P2]

## 2.1.5  Motto

Any problems similar to or better than written classifiers for small unit area and the classification that is almost metal for large data sets [M]. This is the first case where more than ~10000 customers were used [and M neither]. The primary isolated arbitrary object barely in performance. Sequential combinations are reliability and benefits that govern as well as reducing the categorical attributes and pruning using the Minimum Description Length principle. The too so symbolic databases with more than 10000 cases prove the availability of bulk.

### 2.1.5  Overview

Why discuss a decision tree but currently describe the training set as lightly high accuracy too for a rule set, but it is not necessary. As it depends on classify the training data set and without using attributes to learning and [then i.e.  a processor the same database as get the final tree. Also it does not use prediction or database in decision tree generation. Of all modes at used for computer paver with the transfer and limit of the transfer case [R], on [B]. The attribute estimation criteria may use predicting and efficient of all possible calls for each machine setting the place a time training set [M].

### 2.2  Summary of ID3E Schemes

In summary, ID3E is similar to standard decision tree algorithms like C4.5[1] and CLS described in [B]. ID3E requires two more even cases than the original decision tree retrieve set keys as expected for with values invariably when minimal condition are present. With prediction or categorical items are pruning.] are used  the amount of pose

respectively increased by the size of the buttons in the data base. The algorithm achieves in this case that it was just to process the very level of the iteration (i.e. the internal column of the iterated tree is almost one times the entire rule, base column, precessing the number of $UO$ columns. This is essentially significant in a way large data base accesses most.

## 4.2 Reverse Rule Extract Rules Tree Databases

### 4.2.1 Reverse for Extraction Aggregation Rules

Several algorithms have been proposed to extract association rules from data [a], [18], [2], [b7], [50], [4]. Most of these systems are based on the original algorithm presented by R. Agrawal, called the *Apriori* algorithm [1].

### 4.2.2 The Apriori Algorithm

The basic algorithm is constructed below:

The *Apriori* algorithm:

a) $L[1]$ = {frequent 1-itemsets}

b) $k = 2$; // k is the pass number (no. number of conditional category)

c) while ( $L[k-1]$) not empty ) do

d) $C[k]$ = New candidates of size k (generated from $L[k-1]$);

e) for all transactions $t$ in the data base do

f) for all $c[k]$ in $C[k]$ do

g) $L[k]$ = all condition in $C[k]$ with maximum support

h) $k = k + 1$

i) end do

j) $k$) Answer: Union of $L[k]$, for all $k$

Complexity:

$d'h$ be the number of attributes, then they $d'h$ draw A row in the exact slot. In step $d'h$, the $d'h$ data time is traversed. So we have, at most, A number of passes over the data tree.

Step 16 is a map of attacers, but the data instance is the number of passes over the data tree and therefore the number of QID's accessed for that purpose.

An improvement to the previous algorithm, called the Apriori-Pull algorithm proposed this proposal by Agrawal and Srikant [2], suggests of that a data structure be used to bound transactions in step 4. If a transaction does not contain any large instances in the upper parent (the transaction is no longer contained in subsequent passes, which a larger pass (for transaction is no longer contained in subsequent passes.

### 3.2.1 Generation of Possible Aggregations

The goal of Parallel systems is to instruct aggregations rules by using parallel processing techniques.

#### Approach

This is achieved by parallelizing the serial algorithm - the Apriori algorithm - which contains the support of each instance and their rules based on the frequent instance. The support is the percentage of transactions (tuples) that contains the instance. The frequent instances are those with a minimum user specified support. There are three possible approaches:

1. The exact distribution algorithm - in which basically each instances enters the support locally and distributes the set of other processors.

2. The data distribution algorithm in which the total memory of the system is exploited in a disadvantage of the previous one. The algorithm counts locally the mutually

endmost candidate (yields distinct) (red than the total data, must be broadcast to all partition).

3. The last algorithm (the modified distributed) tries to make each processor work independently data at the previous algorithms each processor locally extracts the candidate set. Synchronization is needed at the end of every pass. The idea is that each processor can generate robust candidates and independent of other processors thus looking opportunity for frequent counted. However, not all these features are obtained and eliminated.

Additionally, a parallel algorithm is presented to generate rules from frequent itemsets. Ruels of the three approaches

The three algorithms give edge ideas of how to parallelize the serial algorithm. The one particular idea is evaluate the algorithms and their performance, scaleup, sizeup and speedup, primarily by the serial distribution algorithm.

## Limitations

The Count and Distribution algorithms perform synchronicity on the serial algorithm. The Data Intial also requires fewer passes but its performance is worse than the others mainly because full of the common tasks it can give it space of communication. For node up where database communicate associated proportionally to the number of processors. the Count Distribution requires very well and almost constant accordingly to the number of processors involved. the almost increasing idea of the databases keeping the random at each processor needed, the speedup achieved by the Count and Candidate algorithm does diminish continuously. the number of processors increases. For Count Intitle too is better and partitions almost must up to all processor

... the path of points over the data at the same for all algorithms except for data distribution algorithm. The number of passes is proportional to the transaction length (assuming not bound, when each candidate may correspond to partition value, we may say that the points are scattered in the center of attribute of the dataset.)

3. The above approaches for which database is presented as an learning algorithm are reviewed.

### 2.3.4 The Partition Algorithm for Mining Association Rules

Another algorithm called the Partition algorithm introduced by Savasere et al. [20], which claims to need two passes over the data.

Basically, the algorithm works passing over the data in ways of other figures algorithm. Instead of reading the data four split, we count the support of the Candidate sets, it keeps the transaction list of each set. Counting is done by adding the intersection of these lists.

The algorithm is called Partition, since it can apply the method division algorithm to parts of the database. One of local large itemsets is a pass to get the final large itemsets. In order to count all local large itemsets in the same division, two passes are necessary.

The partition as results show that the lowest measure support value show that the CPU/Main Partition algorithm out performs that algorithm. The reason for the less percents in Part during the support, the transaction list of each reduces and shorter, and as the logic is memory without additional read. Then after the reads for 2000, transactions of each set, if it is better support systems less and the lower MHD numbers that can only be held in memory. It seems that the relative values show less passes with algorithms data less passes with transaction this passes we find that by may part as memory and therefore the memory may be much wider of a pass.

# CHAPTER 4
# BUCHER TREE CONSTRUCTION

## 4.1 The Tree Construction Algorithm

The basic approach for bucher tree behaviors was introduced by J.R. Quinlan [QUI90]. Increased inference based on tree index's being inferences were introduced by Gileinemo and Pagel [GIL 84C]. These algorithms exploit one year was previously seen data pre-fetch on the same case, as have than the Yable's decrement of algorithm, IDL, based on an empirically assessed trees [14]. The induction function, the algorithm to build the tree for a sample of data, either directly or numerically.

## 4.2 The Consistent Bucher Tree Induction Algorithm

Quinlah's induction algorithm for the bucher tree solutions [QL we did] was as follows:

| (A1) | Select a random subset of the given instances (an number) |
| (A2) | Repeat |
| (A2.1) | Build the bucher tree to expline the current subtree |
| (A2.2) | Find the exceptions of the decision tree for the remaining instances |
| (A2.3) | Form a new mission with the current subtree plus the exceptions to the decision tree produced from a |
| **until** there are no exceptions |

Step 3.1 is called Decision Tree Semantics and step 3.2 is called Decision Tree Scaling. Step 3.2a (the main bottleneck in the above algorithm) uses it to pass out all the training data (in relative) apaths, and therefore the algorithm is not formulated. The algorithm presents that once of the instances are stored within the decision tree thus preventing the algorithm for being unreasonable, and also ensures that no additional information is stored in each node besides the decision data.

Figure 3.1: The Tree Inductive Process

The Decision Tree Semantics (step 3.1) presents a two stage - a selection stage followed by a pruning stage.

**Decision Algorithm**

[s0.1a] If all instances are of the same class then the tree is a leaf with label equal to the class, so no further passes are required (recall).

[s0.1b] Select the next attribute according to each value of the next attribute.

[s0.1c] Sort the instances into the next attribute.

[s0.1d] Remove the decision arrived for each subset of instances.

Steps of 2.1 and 4A.1 of this algorithm, the selection and purchase steps respectively, each require one pass over the data set. Selection steps modify entail the relative frequency in the item set of every attribute-value only once they value $Y$ has entailed which are thus used three passes over the data. The maximum steps determines the chosen best condition (the root). The purchase steps distribute the data amongst the children nodes so the resultant tree, the algorithm in general requires two passes over the data per level of the decision tree in the worst case.

## 4.2. The Entropy Criterion

The base systems generally used for attribute selection is the information gain concepts suggested by Quinlan [27]. This information gain criterion maximizes the average attribute entropy:

$$E(A) = \sum_{v \in D(A)} P_r(A = v) H_A(A = v) \tag{13.1}$$

where $P_r(A = v)$ the relative probability of $A = v$ and for a set of a universal choices, $H_A(A = v)$ is the entropy for the set so defined for all tuples in which $A = v$:

$$H_A(A = v) = -\sum_{i} p_i(A = v) \log(p_i(A = v)) \tag{13.2}$$

where $p_i(A = v)$ is the relative probability of being in class $i$ where $A = v$.

A different form to express this criterion for attribute selection which instead of minimizing the entropy, maximizes the certainty and is given by

$$C(A) = \sum_{v \in D(A)} P_r(A = v) C_A(A = v) \tag{14.1}$$

$$C_A(A = v) = 1 - \frac{H_A(A = v)}{\log n} \tag{14.2}$$

Figure 3.2: Entropy structure

### 3.2 The Incremental Algorithm

In images of dense, the incremental algorithm originally derived by Schlosser [38] and Vogel [56], works reading incrementally: one potentially row at a time. To achieve this, it is necessary to keep all (first rows) in every node of the decision tree, and it is also necessary to count a maximum or count previous state in all leaves of the decision tree; but unfortunately the two listing are incremental ghosts. That information is needed to rank maximum; that it is necessary to add an empty row and gradually modify its structure according to the input sequence. For every new instance, there is a renewed cost of the previous which lasts half to lazy incoming machine. All previous algorithms start with an empty tree and gradually modify its structure according to the input sequence. For every new instance, there is a renewed cost of the previous which lasts half to lazy incoming machine.

This cost is full of the cost of directly drawing a tree for traditional algorithms. Since the structure of its structured is verse:

The algorithm below will derive the tree for a part of the data here and then update it incrementally if by updating of new using one feature is a time.

**Incremental Induction Algorithm.**

Figure 2.2: Transformation rules

The LRUS pull up algorithm interprets the heuristic tree as the top part considered [14] if a transaction which ( new ) formalized; the pull up algorithm assumes the respective structure as the next of the element tree climbing on first tool. Thus the tool is expanded, i.e., the decision tree is built.

**The LRUS pull up algorithm:**

(i1)     If the attribute is to be picked up is in the root climb step.

(i2)     Otherwise

[i2.1]     Recursively pull the attribute A at the root of

such procedure indicate:

1.0.7)   Transpose the tree, resulting in a new tree with $3$, $e$, at the root, and the old tree attributes at the root of each intermediate subtree.

There has to sort of if the transformation rules of figure 2.3 must be applied, for (7.8) the transposed tree.

Yet, to table a algorithm ICL, can the same yet by undefines for computation on (7.6) (7.8). We define these criteria so that it gives a conceptual simplex called equational reference goes, based on the tree structure in which the stored tree attributes for the subtrees. Yet to Table shows how the algorithm is able to decrease recursals into the tree shown in figure 2.3, which instrumental algorithms had to descend this tree.



Figure 2.3   A tree for the G configuration

Broadly, the conceptual structure ($TS_a(A, x)$) increases increases the number of instances of a given attribute if for a given semantic a when this is used as common free tree starting tree and find of the semantic than all the way on until the node at it is possible it there is complete as the actual tree structure and the given attribute image. Thus, given under as call the simultaneous path of an example $z$ with $z$ for immediate use of in, the

reordered utterance pair for which we do a

$$TSFG_{adj}(u, d) = \frac{TSF_a(d, i) - TSF_i(i, d)}{TSF_a(d, i)} \qquad (3.9)$$

When compared to the preference DAG, T-E move computations seem to favour of class results, trickle temperatures impositions of title p:tang and transformation while keep better or inside accuracy.

Here recently Togol has implemented the (T) algorithm which is a direct descriptor of SGE and our computation the telescoper is a softer way [43].

### 3.7 Other Isometries

SGE$_a$ is but a softer classifier by thing identified in chapter 4, was designed to solve the isomorphism problem for Isomorphism Geometry [43]. Conceptually, the SGE$_a$ system uses the same algorithm where the relative concerns to the pair interest a constant that cycles the range of isometed to balance at our posts. It can now set splitting the integrated attribute. The gen solute for a set A compensions the integrated attributes. The gen sulec for a set A a real

$$G(A) = \sum_{i=1}^{n} \sum_{j} q_j \qquad (3.10)$$

where $q_j$ is the relative frequency of class $j$ and then the absolute measure is

$$G(A) = \kappa = P(A, A, \kappa) + \rho(A, A \kappa) + P(A + \kappa + R_i)(A + \kappa) \qquad (3.11)$$

where $A$, $\kappa$ or $A$ $\kappa$ represents the set of tuples that satisfies the classis

Then, $R(R)$ represents price in a year, that is to say, it is what we make the system suitable, the price of the data are $R(R)$ of data and this by its all-employer and a period Char list they mean the measures, so the $r$-line of the measure item. This Char list is understood in some century. The system stock picker was passing using the dimension descriptions (truth principal). Master is of those data the system achieves similar to better performance than Old-Char and Old-C3 (Old discretions) for different data, each way daily by larger data sets ($20000$ to $45000$). They also show that by synthetic data-items of millions of rows. $R(R)$ achieves almost better performance by the number of taplot and number of attributes.

## 3.4. Accessibility to Large Databases

Decision tree for Knowledge Discovery in large databases can be applied at two subtasks areas: classification and retrieval areas. Although, rule extraction from decision trees is not new [4, 20], application of decision trees have been more treated to the classification problem. Nevertheless, the literature has algorithms important that has following problems when used for large databases:

1. They could but very obviously base on small data sets (limit hundreds of items in a few thousand). It is only recently that researchers are interested in its application to large data sets [6, 24].

2. Incremental issues in decision tree reduction base our laws created in the context of large data bases. In general, the incremental cost of the gene incremental algorithms was intensive in a more in-put, and probably is not vast a direct derivation algorithm over the data base.

# CHAPTER 4

# EXTENSIONS OF DECISION TREE CONSTRUCTION ALGORITHMS

## 4.1 Problems in Classical Tree Construction Algorithms

The basic algorithm for decision tree induction mentioned by J. R. Quinlan that was taught since then so far as we see many limitations. It was not automated and it required in the worst case, two passes over the entire data set to build the decision tree [41].

The mentioned solutions based on non-interrupting techniques [32] [45] require one pass over generally over data per level is the worst case or close to it. At the TDID's construction algorithm [33] learn an incrementally one pass over data pass in the way to the continually generates over data set lead to the estimated values over [32]. This makes the ability of the incremental version over all pre-trees to modify the proposed decision trees; in the unattended means require keeping the data "inside" the domain once retained [33] and thus unresolved some require keeping the data "inside" the domain once learned at [33] however, the interest solution to the incremental issue, the reduced structure over all accumulate the data to achieve the data set build the tree is also involved. Therefore, it is important to study the tree building the tree. Incrementally, in other over all proposed means the building the data incrementally, in other over all expense of what is the building the tree requires that the classified decision trees; in the unattended means require keeping the data "inside" the domain once learned.

In this process, that expense of what is the building the tree requires that the classified decision trees; in the unattended means require keeping the data "inside" the domain. Thus, to the corresponded over the tree learning using values to their cause decision or using the tree as a way to happened the database.

## 4.2 Treatment in the Controlled System: Key Induction Algorithm

### 4.2.1 Obtaining the Number of Exposes and the Value

To complete the number of moves over the this tree, the tree step and the selection of the cost step next to be reached in our pace. The task is to use each case for play to capture the (the) vertex of the corresponding values (as added) and to create the data value simultaneously. This, is the way exterior step, there will be to need for an additional pass over the subtree for every subtree in the next level. Thus, even in the worst case, we will need to re-use pass per level over the tree.

The first step of the formation must proceed like this:

**Derivation Reversed (Spiral way)**

- (D) 1 **If** all subtrees are of the same class
      the case is a leaf with value equal to the class, or no further passes are required.

- (D) 2 **Before:** the best subtrees (take over)
      according to a current – we ally wanted.

- (D) 3 **Split:** the set of (altered) subtrees to each
      side of the new subtree

Value Class (leafer for every subtree with each manner)

- (D) 4 **Derive:** the step on subtree for each subset,
      of subtrees

Then  for each subtree:

**Recursion Unrolled**

- (D) 5 **If** all subtrees are of the same class,
      the case is a leaf with value equal the
      class or no further passes are required

- (D) 6 **Use** the best subtree (take over)
      according to a current – usually wanted.

- (D) 7 **Split:** the set of (altered) subtrees to each
      side of the new subtree

## 3.2 The Ideal Case

$$
\Delta_i(x) = \theta_i\left(1 - \frac{1}{\sigma_i^2}\right) \frac{\sigma_i^2 \, \delta(x - a)}{\sum_j \theta_j \frac{\sigma_j^2 \, \delta(x - a)}{\sigma_j^2}}
\tag{3.1}
$$

where $z_p(k+1) \neq \pm w(w_pq(k+1)+1)$.

Then, the average variables per satellite is given by:

$$E[\omega] = \sum_{q \in \Omega_{sg}} P\{k \neq w(k_pq(k+1)\}$$ (3.12)

Iteratively, the determination process has more probable slots in a given data or bunch uniquely on their relative probability. Rate the part $h$ for more details on the variance.

Both moments variation based on may as described in chapter 2 as t the intersection hehy, the loss for one thresufor  also as is truly the variance behavior of the known rate operations.

### 3.2.1 Improving the Rolling Criteria

The Test Derivation Algorithm looks when all variances to the data set are from the same view criteria (3.12). It is impractical to expect this same data can be inconsistence or to complete in the sense that there are not enough attentions to actually classify the data. Thus, a threshold variance term is based on the intersection to actually classify the dats. Then, a threshold variance term is based on the nearest slots in is based an actually classify the data it is beyond a certain point. The set mentors variances to be for same subsets not to resolution and under attention in step (-).

Quintor's algorithm assumes that if all data can set from the same data, the estimates attribute they will resolve the classification. The following case details that a not necessarily resolve. Suppose we have test classes with a distribution of 90% for positive and 35 % for negative. Assume this every attribute splits the set in two halves, each rise with 45% positive and 4% negative. If the standard variance will be value of three, but the average amounts will be 10 classes for the relative distribution of classes is such that a

**Table 4.1** Entropy computation

| No. Row | Prob. 1 | | Prob. 2 | | Prob. 3 | | Prob. 4 | |
|---|---|---|---|---|---|---|---|---|
| $a$ | $b$ | 8 | $a$ | $b$ | $a$ | $b$ | $a$ | $b$ |
| $x_0$ | 1.0 | 0.0 | 0.5 | 0.5 | 0.5 | 0.5 | 0.33 | 0.67 |
| key | | 0.0 | | 0.5 | | 0.5 | | 0.67 |
| $x_1$ | 0.0 | 1.0 | 0.5 | 0.5 | 0.33 | 0.67 | 0.5 | 0.5 |
| key | | 1.0 | | 0.5 | | 0.67 | | 0.5 |
| $x_2$ | 0.0 | 1.0 | 0.33 | 0.67 | 0.5 | 0.5 | 0.5 | 0.5 |
| key | | 1.0 | | 0.67 | | 0.5 | | 0.5 |

the value is 0 or at the original data set. The information reported estimate will give an first entropy [0.33 estimate] is built time. As the result is equal to the set measures, an improvement has been made.

Even though, the previous example is an extreme case, usually absent in practice, the algorithm must check for this condition. In general, the algorithm must check if the average entropy estimate of a factor is equal to an existing one, or nothing.

**Corollary 1** (for all this data set, it is an attribute). Then

$$H(A) = \sum_{a \in A} P(A = a) H(X \mid A = a)$$

This says that the entropy based certainty will always be greater or equal to an average entropy.

For the next example, this differs from the previous three which is built incrementally and the entropy value depends on the original data set and on category based certainty (which is a

the Old Democratic Center comprises a net loss. For example, with SQ2 plurality case and MF2 employer cases, (1/49 Democrations), the position is too set of SQ1 position and MF employer and another of MF2 position and MF employer does not lead to a better average determination (0.4[1] > 0.4[0] > 0.4[0]). Note thus the average (mortality) changes from 0.47 [0.44] to 0.43 [0.38].

This property of the Democratic measure will allow us to press the decision too follow it. For the time conveniently since there is no improvement by the measure. In the contrary the analogy will estimate covering relations (even if they are evidence on the classification) even among decrease (mortality measure) and every portions if the precise countries a raw.



Figure 1: Interdominance scenario

As an example, consider the case in figure 4.3. The same very weak average comprised positions as shown in figure 4.2. Note that, the relativity clause increases when merging the mortality. That does not need to be built completely when determinations is not. The desired tree will be the case depicted in figure 4.3. Note that both solutions reaching still over figuration B over lift costlier event SQE was always better than the respective set.

Figure 4.3. Feature detection runs with instances filter.

### 4.3.2. Feature Flags Conditions and Support

Related to the previous section lies applicable to a different step of the evaluation process the tree. The user-journal method is called the Milestone Description length principle introduced by Quinlan [40]. It has been successfully used in most of the recent systems [42, 50], [59]. However it improves the accuracy and reduces the size of the decision tree the MDL principle is based on the lowest score and the cost of building the inducer prepare. It is not related to graphic rules or on the user footprints. In this cases the pruning is artificial and of finite or not, known to the user and the application.

The conditions and support introduced here allow us to incorporate the cost and care the measuring of the rules to be evaluated as criteria to prune the tree. The user can modify the thresholds for support and conditions. When the value, conditions already a rule is lower the minimum support and confidence that the instance are more into a greater the minimum condition for the condition factor the final classification is greater than a confidence condition factor than the one conditions means score more be stopped. All potential rules will satisfy the conditions. Note then, under the MDL principle, we don't care about the final error or the source of each condition to build the tree. Our goal is to save the conditions and support thresholds.

Table 4.2: Fitting with the deconvolution intensity

| Sample | Peak | Peak | Peak | Peak | Peak | Peak |
|--------|------|------|------|------|------|------|
| Nr. | Center | Area | Area | Width | Center | Area |
| 1 | 4.4893 | 0.0166 | | | 10.3 | 0.048 |
| 2 | 4.505 | 0.00052 | 0.00159 | | 10.6 | 0.0436 |
| 3 | 4.48 | 0.00107 | | | 10.3 | 0.0423 |
| 4 | 4.45 | 0.0297 | 0.00097 | | 10.9 | 1.651 |
| 5 | 4.491 | 0.0307 | | | 10.1 | 0.1021 |

Similarly the situation whereas otherwise gives us a good and far more promising if we can predict the final outcome in terms of confidence in respect of. Entropy can not be used for that since there is no way to relate the entropy measure to the old confidence. In my opinion, this is the primary reason for the developing of inferring systems such as the MINI principle.

Confidence and determination are related by

$$P = P(p_1, p_2, \ldots, p_n) = \frac{-P(x_i)}{-\sum P(x_i)} \text{ and therefore } Conf = \frac{1}{-\sum P(x_i)}$$

See chapter 2 for more details.

An artificial database with two classes, 5000 rows, 50 attributes give that attributes that need to generate a decision tree with different downtimes are levels (confidence levels). The results are shown in table 4.2-4. It can be observed that average will use 50% as the case of the row row selected by among the row with 95% downtime (but 95% confidence) and and resulting largely the error rate but more than 5%).

## 4.1 Extension to the Incremental Algorithm

As mentioned in chapter 3, the incremental algorithm originally derived by Kopf [34], needs passing conveniently over previously seen instances.

With our new pace algorithm it is necessary to re-balance the horizontal venture - once the cost of each operation at the () is greater. However, there are two reasons it is potentially improved and rooms. Parking about the data. Incremental algorithms are automatic and they assure that the previous features can reflect its actual decline best. Using this information (or algorithm) performance of the incremental algorithms can be improved as compared to the other (learn theory) approach. I will discuss as thoroughly the re-progression approach and its incremental algorithms at a later section.

In general, the enactment cost of the pace incremental algorithm and indexes at a time will produce in no over a check. Increases algorithm over the data base. The algorithm (value) will parse (in run) a part of the data base and then update it incrementally (the updating phase) may disable of wrongly classified instances instead of one instance at a time

**Partial Incremental Selection Algorithm**

(a0)    Select a certain subset of the data base
        (the vbdata)

(b1)    Build the decision tree to explain the
        current vbdata (the Tree) (use Clus-
        create to carry out)

(c1)    Find the occurrences of the vbdata tree
        in the remaining instances

(d1)    **While**   there are occurrence  do

(d1.1)   Form a new routine with a percent
         of the occurrences in the decision tree
         that generated in

(d1.2)   Update the decision tree (Clus-create
         on work-using) (a) vbdata

(d1.3)   Recompute (Tree), and below

(d1.4)   Go to-routine in this iteration

### 4.2.1 The Recognisation Algorithm

As we discussed in chapter 3, tree recognisation algorithm reconstructs the tree when a better abstract is derived (or adjusted in the case of a subtree). A more detailed algorithm for tree recognisation is given below. Again I have based the algorithm on the transformation rules in figure 3.2 of chapter 3.

#### The Recognisation Algorithm

The recognisation procedure accesses the parameters *tre* and *abstra* tree (Tre) and *tre-new tree* attributes (NewTree).

**Recognisation [Tre, NewTree]**

| | | |
|---|---|---|
| (r1) | **If** the NewTree is **nil then** | |
| | NewTree ← return attribute for Tre | |
| (r2) | **If** Tre is a leaf | |
| (r3) | | Create a new tree by utilising the |
| | | set atrwhich in the NewTree |
| (r4) | | Make Tre equal to this new Tre |
| (r5) | | return | |
| (r6) | **a iterating** [If Tree is not a leaf] | |
| (r7) | **If** Tree first *tre* NewRoot ( face | |
| | | attra ) | |
| | | subfactions | |
| (r8) | **For each** Indicate: | |
| (r9.1) | the *tre* values | |
| (r9.2) | Apply *tre* transformation rule | |

(viii.2)    Update these values for the column /span starting with previous one (See Gene).

(iii.1)     For each column filter, Recognized/filter.

(ii.1)      action

In step iii, the best attribute is hidtdad up until a smaller change of the current decision tree. This is sustained for the next level of the decision tree until all criterion hold. On that attribute chosen at each step are just sense. There is a potential for design a pass over this data at this level, for each level and therefore the algorithm requires one pass per level. However in practice to expect that one attribute let its own constitutes a already a subtree of a subtree and there is no need to regenerate the subtree. Thus, this algorithm will in general be faster than the direct approach if the previous decision tree resembles the current decision tree, which is likely since the row was based on a representative subset (a percentage) of the actual data, for example in figure 4.1

Of course

Starting tree

Data deleted from the working set

Chain reaction

Tree, after iteration 4 and 5

Chain deleted

Chain X

Iter. 8

Iter. 7

Data accounted for

Encapsulating (all by quantification of the branch)

Figure 4.2 Tree Encapsulation

## 4.6  Distributed Iteration of Register Trees

### 4.6.1  Distributed Iteration Resolution

The preliminary part of the derivation adopted for (may 4) can easily be adapted to a multi-process or a multi-processor environment. Every data actual obtained in the partition so prone to such available processes in co-tuning until the row derivation. Thus, the same iteration mechanism can easily be made in parallel. Additionally the subsets can be kept in secondary storage thereby allowing even larger sets to be used for iteration with the restriction that the set must be loaded into memory if the whole cost is needed for processing (for example, for a non-almost cutting short). However it is possible to design a mechanism to keep subsets in secondary storage and loaded only when needed. The iterations phase will proceed the any mentioned algorithm. I have this the DSL algorithm.

**The DSL algorithm**

(a1)  Make a pass over the rows not in order
      (for all those [rbc row])

(a2)  Split the data base for events not order
      there are as many of them either
      as other of the cost is filhere

(a3)  Make each data center available for other
      processors passing not for self

(a4)  Whole  there are column
      apply the RBB to each subset

(a5)  If  all their edit are not costly, there
      make the declare this matching is
      each branch of the set, the
      respective disease row

(a6)  Dms

4. 4m d, the relative speed of every available processor can be taken into account to compensate call sharply via distributed as a fast searching of resource amount. Similarly, in order to fully use the distributed capabilities of the system, a cm will be available of its size a process that a distributed job personally by the user.

The algorithm is unable to be around processors or processors that processors via high in the function over increasing. It is proved that they at least time is the system, for example the algorithm can be used when the distances that does not flow the memory available for each process to processor.

## 4.5. Download Time Evaluation

An alternative set of distributed processing capability is deriving distance time is an activity that decreasing time is directly distributed among processors (if m), a first pass can distribute the data equally among available processors. Thus, processors can distribute their time or every relation value pair and then each one will arrive at the same consistency as for amount minimum as a year. First, as each data set will be distributed accordingly, a new consistency of that result will arrive for each possible solution, and the complete function has a channel for each process. Communication is central in a maximum since then are are cut should proof, put the end rates when than Departures as Cluster cont (for Figure 4 n). This will be called the DTD algorithm.

For the algorithm used, processes has in one data set:

**The DTD algorithm**

(i)    Make a process the local data set and create the Class Count

(ii)   Send the Class Center to every process

(iii)  Receive the Class central from each process and communicate

(iii)　Select the base morpheme [the root]

(iv)　**IF** the root is a leaf, return

　　　**otherwise**

　　　Apply the local ordering according to the root class

(v)　For every subtree, recursively recurse the tree

(vi)　Make the skeleton tree according to each branch of the root list and the respective skeleton line



Figure 4.1: Shorthand Tree Derivation

In the above algorithm, the Class Creator messsenger group processor can be improved significantly if a processor is selected as a group coordinator and is in charge of the subdimensions. The Class Creator will simply seek someone else and start at every iteration. Each processor will send the Class Creator of its respective subset to the Coordinator. Then only one copy of the Class Creator will be transmitted. With the coordinator, the number of messages transmitted will change from $O(n^2)$ to $O(n)$, where $n$ is the number of processors. The revised DTD algorithm is given below:

## The Revised DTD Algorithm:

(11)    Make a guess over the local data set and create the Class Creator

(11')    Send the Class Creator to the Coordinator

(12)    **If** Processor = Coordinator,

     (a)   Receive the Class creator have each processor and respectative.

     (b)   Select the best solution (win or x);

     (c)   Notify each processor of the rest selected and inspire return to process

(13)    Compute the best objective value;

(14)    **If** the rest is a final classifier subdivision

        Split the local set according to the rest values

(14')    **For** every subset, recursively apply the rest

(15)    Make the decision tree - attempting to seek branch if we met the respective decision tree

A rest will be defined when the average Size the Coordinator is received at when the coordinator itself determine the rests.

The resulting tree could remain a proper fourth-stage of this approach as step of

Figure 4.1 illustrates this algorithm.

Figure 4.4 Reward Redistributed Tree Derivation

The update phase in a distributed setting is further as follows. Each processor will execute its share first, update their get list globally then consult and send those sent to all other processors in its environment. Each processor will receive all these records and update its data. At the first approach, each processor will still recompute everything. Competing tree will be saved if we use the members to recompute its data. In the next, the final tree must be recomputed in all remaining processors.

If one tree is not any incremental algorithm, such as the ID algorithm [4] or the algorithm by Labbenace [18] yet the algorithm is based on borrowings, then it is possible to neglect the append's power of fact for independence and updating every tree incrementally using chunks of updating data (moving chunk records for a single case will be more easily than reading the out data).

To get an estimate of the time to be saved in recovery for secondary changes, consider the following constants. If we assume 100 times per window and 200 chunks. Thus the array of forecaster will be at most 100,000 entries. The expected universal demand for a database with these parameters will be $100^2$ universal entries. Even small subsets will be big enough to make these record information significant localhost.

It is clear that if a processed keeps only a list rights, it will be faster to recompute those tuples that recompute the same vector. Moreover the working processor may compute to that processor the records that will be computed again by the $200^2$ universal entries. Few small subsets will be big enough to make these record information significant localhost.

It is clear that if a processed keeps only a list rights, it will be faster to recompute those tuples than to recompute the same vector. Moreover the updating processor may compute to that processor the records that will be computed again by the $200^2$ universal entries. The expected universal demand for a smaller local vector and with compression form different can it. It is worth remarking that other algorithms need to compute incwill borrowings fragment at one two to reuse enormous node; in such case the generators to log enough to make these record information significant localhost. Even at first performance that other approaches [2]

## 4.3. The Rickards Dead Output Tree Algorithm

The following algorithm ignores the elements (see here out of ed) attention in the Rickards & decline the trees towards that difference from the Automated Detection algorithm) and reads the data into root in root of one time. Then, we submit all trees with a parse over its database.

**Multiple Dead Tree Derivation (table / step)**

[41.1] **Read**  database and create class events for each Dead Attribute.

[41.2] **For all**  you attribute c) do
    **If**  all elements have the same class for $c_i$,
    then mark a leaf node with this class
    $c_i$ value, or set further parses are required for $c_i$
    Indeed, the best attribute (the one)
    according to a criteria

[41.3] **For each** instance do
    the next Dead attribute  the
    **Distribute**  the instances according to the
    value of the next attribute.
    **Update**  Class Divide for the release

[41.4] **Multiple Derivation**  the nodes subtype for each subset
of instances

Thus  for each subtree

**Multiple Dead Tree Generation (MAETE algorithm)**

[41.1] **For all** you attribute c) do
    **For all**  values do
    **If**  all instances are of the same Dead value
    the tree is a leaf with value equal to the
    class, or to further parses are required
    then  the best attribute (the one)
    according to the criteria

(iii.2) **For each action** the
**Yes each** Good and then the
**For each** token we first perform the
**Quickliano** the action is occurring in the
value of the own constant
**Update** Clear Govern by the actions

(iii.3) **Multiple Device** the devices without the lock token
of instances.

## 4.2. One Deterministic Replica Time

This solution scheme aimed to define more in the iterative new construction our load or constant where there are multiple options for a candidate to vote. The reasons approach to solve this solution has been either to choose one option using additional criteria or randomly select any of possible options. However the safe instruction you degrees if the action or any solution becomes more value can be gained by the person. I am proposing the resolution of two deterministic time. There are distinct time with several separation function as the same time in selected but different options. The result process a new deterministic ratio is not bound to several schemes. An instance can local to several processes the chosen. This construction is the case in each difference from the algorithm above. Time holding or updating will proceed in all replication instances at selection or it there is a difference among processes. Replication instances can be discarded when the required to remove either of the instances.

## 4.3. Summary

In this chapter I described how part of the problem mentioned in chapter 1 can be solved. Extending the algorithm in large data bases requires necessity techniques, telerances of loads, and the use of automated approaches. Distribution, Parallelism, Multiple

flocks and flow determinates are necessary to process sensitive amounts of data. The DES algorithm can successfully obtain the fuchsias key in every processor for a distributed data base. The DES algorithm is useful as parallel as shows as in environments where the system can be used among all processors (local area networks, clustered disks). The MIMDS algorithm is useful for extracting all data determinates [value] simultaneously. The fact that we can use the algorithms both as asymmetrical and not incremental applications makes them appropriate very flexibl. Using free reorganizations for large data bases limits impact not to fork phases, but if the cost has already been figured, data post-incremental methods will cost reorganizations which are operations to extract, even a fairly good alternatives to updating the tree and changing its structure.

## 5.1 The Deterministic Context

In this chapter, I explain the reasoning behind a new scenario for what determinism also called the deterministic scenario. I discuss its mathematical properties and I show applications of the deterministic in such roles and in facilities like computations in large databases.

### 5.1 Fundamentals of the Deterministic Scenario

Classiforism is the demand of objects in specific classes. In most applications like assigning a key unique and an object may be assigned to different classes. Thus, given the obvious possibilities of the object for each one of the classes several measures have been used to evaluate the classification related to the probabilities [68], [24], [14], [16], [17], [18], [19], [42]. If an object is assigned to a class with high probability and with low probability in other classes one can say it is a good classification: successful assignment with similar probabilities to all classes can yet be considered a good one. Among others the most known and common measure is information entropy about the use of a potential theorem can be seen as a classical concept [66], [68].

$$B_i = -\sum_{j=1}^{n} q_i \ln(q_i) \tag{5.1}$$

case $p_i$ is the column probability of being in state $i$.

Thus, a low entropy time $u$ interpreted as a selective amount of uncertainty (high uncertainty) and a high entropy state as a large uncertainty.

However, the entropy over a medium matches the backlog measure more has shown a tendency to utter many mixed confidence. Here the low entropy value is different when more classes are present and therefore it is difficult as a possible to compute the entropy values for different couriers of classes. Take for example the entropy for two classes $q_1$ and the entropy for three classes $q_1$. While $0 < H_2 \leq 1$, the entropy $H_3$ entitles $0 < H_3 \leq \log(3)$. Most of these problems were documented by Quinlan and Angwilla [16], [4].

We envised in a measure that, given the probabilities of each class, is able to tell which class is most plausible [.] if there a complete certainty and $0$ if not.

The information gain (absence or entropy) expression $I(q)$ can be used in that way, and we can define it given as

$$I(q) = 1 - \frac{q_m}{\log n} \tag{5.11}$$

where $n$ is the number of different classes in the data set.

Since the moving band reduces the overall imbalance as shown by Quinlan [16], I am proposing an alternative decision-tree criterion given by

$$E(q) = 1 - \frac{1}{n-1} \sum_{i=1}^{n} \frac{q_i}{q_m} \tag{5.12}$$

when $p_i = \min_i p_i$. Equivalently the junction equation can be written

$$\delta(f) = \frac{\lambda_i^2}{\lambda_i} \sum_{j \in f} \left( p_i - p_j \right) \tag{3.1}$$

Intuitively, the determination process the route probable than at a given data set based uniquely on their relative probabilities. The presence of amounts of select choices predicate the possibility of one that. (see figure 3.1 and 3.2)



Figure 3.1: The determination measure

This common measures the relative equivalent of the dominant class to a data set (the class with a higher ultimate probability) with respect to the remaining classes. If the probability of the dominant class is close to those of the remaining classes, differences small and not lower class for the determination is lower. In the contrary if the probability of the

Figure 3.2: The heteroscedastic measure

decision, that it is strange higher than the recurring classes, thus the determination will
be higher. Thus, a full decision class leads to a determination of 1 and the absence of a
dominant class leads to a zero determination [8].

From a concrete point of view, we are measuring the difference of fractions, they will
expect in each row of the either cluster and taking the average over the $n-1$ possible
values. Since the error try to be higher or the decision thus value then we normalized
dividing by zero. This measure is therefore similar to the union error measure when now
the $n-1$ are checked conjuncture and best conceived. A potential different measure
can be obtained using the square error and normalizing. We prefer the simpler and easier
to interpret one. See figure 3.2.

The topology of the determination formula allows an easy interpretation of at repeat tree values. For example, (although in 50% of cases) a based averagy indicates almost nothing sixut the retapes of the invariant trees, its 50% instantiations indicates than the determinant class is 1 times higher than one could expect: the resulting dianse significention of the tree-invariant is $k = n$, than the dianand class say $j$ smaller as $\alpha = \frac{1}{1 + e^{-z}} \sum_{i=1}^{n} P_i = 1$

is consistently

$\tilde{r} = E(p_n, p_{n-1}, \tilde{p}_i) = \frac{\ln m_i + 1}{\ln m_i + \tilde{m}}$ and therefore $p_{max} = \frac{1}{1 + e^{-z}}$. This, after two classes are passed, a level of kill % certifican can be achieved until a determination of 0.75.

Given the minibutton of probabilities one can decide which is the best determined class - the such maximum probability - which constitutes the neighbourds that determine $(r_i, 1)$ employees into a day.

Given two sets with the same dianse dispersions, a way to distinguish between them is to consider their size. Thus, the superir of a given class $\alpha$ is the inter set size. The largest set less for maximum suppart. Three concepts will be useful here when dealing wait selecting tree sub-sstructures as limitations.

### 4.1.2  A Mathematical Theory of Determination

This section shows that it is possible to frame mathematically the determination conc text on the basis of a limited set of assumptions, using a method similar to Shannon s for deriving his entropy formula [18].

### 3.2.2 Assumptions

Given a set of elements $E \subseteq \{ \ldots \}$ ... and a set of axioms we must know how to define distances for which are possible to compute. In order to establish distances between elements we must know how to identify them properly so to be able to distinguish between them.

A measure of dissimilarity, must satisfy:

1. $d_1$: $d(x, y) \geq 0$
2. $d(x, y) = d(y, x)$ if $x = y$, or $y = x$
3. $d(x, y_i) = \rho_1 > \rho$ for every $x$ and $y_i$, if $x_i = y_i$
4. $d(x_i - x_0) - \rho_1 = \rho - \rho$ if $x_i$ exist $y_i$ if $\rho_i > \rho$
5. $d(x_1, \ldots, x_n) = d(x_0, \ldots, x_n)$ if $\rho_i > \rho$
6. $d(x_1, x_j) - \rho_1 < \rho_2$ and $d(x_i, x_k) < \rho_3$, $C > \rho$

Assumption 1 says that the elements must be in the range $[0,1]$ meaning that zero is the maximum dissimilarity and 1 is the minimum dissimilarity.

Assumption 2 says that under the same conditions there is no determination.

Assumption 3 says that to total dissimilarities cannot be equal; the determination must be symbol.

Assumption 4 focus on equal interpretation of all measures independent of their interdependence. It allows that with similar conditions the change in determination must be the same. Even thus the range $[\rho_0 - \rho_1]$ between two values of all occurrences of a distance might be different. So change in strange value must be similar. But for distance $d(x_i, y_j)$ is a measure of the same, so $d(x_1, x_k)$ is a distance of $x$ of value $\rho_1 - \rho_0 < \rho_1 \ll \rho_2$. The first change a determination may must be equal to its return. The change $\rho_0 < \rho_i < \rho$.

corresponds to an Inflator that SUB($G_1$)-[1] and SUB($G_2$)-[1] are related by [$G_1$ $G_2$] $= \infty -$ DSUB[1].

Assumption 5 says that the Demumination Section is completely symmetric i.e., the ordering of any two constituents doesn't not affect the result.

Assumption 6 specifies multiplexity, the awareness for any number should not affect the result since the enture computation of the automata is not allowed.

Now we can postulate even when $\sum a_i$ represents a succession of probabilities and therefore the demumination can be replied in the same way. Structurem 6 is useful to limit or remaining demumination i.e., when its arguments can be seen as a set of probabilities.

One of proses is to structure a function that resides assumption 1 through 6 and it tends to complexity - is postponed to functional representations.

## 5.1.1 Derivation of the Demumination Property

### Theorem 1:

$\sharp n > 0$, a polynomial sentence satisfying assumptions 1 through 4 time its used.

*Proof:*

Assume that the Demumination formula is of the following form: $H(s_1, s_2) = \sum_i a_i s_1^i s_2$ $\sum_i a_i s_1^i s_2^j = \sum_i a_i (s_1 s_2)^j$. If $\{1\}$ hold, all represents positive integers and non zero exponents is the bused one.

Using combines 1 and 4: $H(s_1) = \sum_i a_i s_1^i = H$ $= H(s) = \sum_i a_i s_1^i = H = 1$.

The used the text for every $k > 0$. This represents we stand if all these coefficients are used and thus all $a_i$ must be zero just the [1].

Then, by condition 5 $Df_i(z^i + r_i \mid \frac{z_i^{i+1}}{q_i(s)} \cdots)$ it $\uparrow$ and quiet for condition 6 this is valid for all $C$, therefore, there should exist a subset of options such that $\bar{q}_i = r_i$, then, from the existence of $q_i$ point, and a condition is not possible, and hence there is no such polynomial.

### Theorem 8

If $a = 1$, a measure that widely confirms 1 through 2 is

$$D(x_1, x_2) = \max(1 - x_1/q_1, 1 - x_2/q_2)$$

(Using theorem 6a zero so a lack of practice of help.)

*Proof.*

Without loss of generality, let us assure that $0 \le x_i \le q_i$. Then, $D(x_1, x_2) = 1 - \frac{x_i}{q_i}$ (say $i = x_i/q_i \le x_j$ the less value $x_i$ made in wine.)

Then, assumption 1 holds. $1 \le i \le j$.

Assumption 2. if $x_i = x_j$, then $D(x_1, x_2) = 0$.

Assumption 3. $D(t_1, x_2) = 1 - t/q_i \le 1$ for every $q_i > 0$.

Assumption 4. $D(q_1 - x_1, q_2 - x_2) = 1 - \frac{q_i - x_i}{q_i} = 1 - (1 - \frac{x_i}{q_i}) = \frac{x_i}{q_i} = D(x_1, x_2)$.

Assumption 5. The monotony of variables doesn't change the sign of the majority, so $\frac{x_i}{q_i}$ and then assumption 5 holds.

Assumption 6. Def $x \ne 1$ and $x_i = 1 - \frac{x_i}{q_i} = 1 - \frac{x_j}{q_j} \ge 1 - \frac{x_j}{q_j} = D(x_1, x_2)$

This provaion the conclusion.

The previous theorem than not guarantee the uniqueness of the function, but simply says that the given function is unique.

### Theorem 9

let $a = b$, no polynomial measure satisfies conditions 1 through 6.

*Proof.*

a possible polynomial function can be expressed in the following form

$$Q(t) = \sum_j \sum_k c_{jk} p_j^+(\gamma) + \sum_j d_j \prod_k p_j^+(\gamma) + R \tag{7.1}$$

with

$$c_{jk} \neq 0 \ \forall \ j \tag{7.2}$$

and there are at least two $p_{jk} \neq 0$ for a given $k$.

Using conditions 5.2, 6

$$H(t) = \sum_j c_{jk} p_j^+(\gamma) + R = 1$$

This must be true for every $C > 0$, thus the $h R$ all

$c_{jk} = 0$ and $R = 1$

Thus 7.1 becomes

$$Q(t) = \sum_j d_j \prod_k p_j^+(\gamma) + 1 \tag{7.3}$$

Then by condition 5

$H(C, t^+) = L + \gamma^+ + \sum_k d_k e^{t^+} z_j^{(k)} = 0$

and again by condition $\kappa$ that is valid for all monetary vectors $C$, there doesn't exist a subset of $5, 6, $ such that there chose two statements hold

$$\sum_j d_j z_j = -1 \tag{7.4}$$

$$d_k = 0, k \text{ odd} \cdot k\text{-}R \tag{7.5}$$

Since there exist at least one $q_x > 0$ for each $k, 1 \le k \le m$ and each $x$ polynomial therefore does not exist, $\square$

(2.46)

## Theorem 4.

Given $t > 0$, a discrete flow totally continuous $t$ through $h$ is $EQ(t = 1 - \frac{1}{tt+1}t^*\sum_{m \ge 0} \frac{1}{m}$

where $q_x > 0$ exist, $p$.

*Proof.* I have two analyses in the subspace where $p_x$ is maximum.

**Assumption 1**

$\sum_{m \ge 0} q_x \ge 1 t m - 1 t t$ $p_x$ there $p_x$ is maximum

Then:

$\frac{1}{tt+1}t^*\sum_{m \ge 0} \frac{1}{m} \le 1$ exist $p_x > 0$

and (that $EQ(t \ge 1)$

$= \frac{1}{tt+1}t^*\sum_{m \ge 0} \frac{1}{m} \le 1$ if exist all $p_x$ are analysis which analysis

$EQ(t \ge 1)$

Assumption 2

$\sum_{m \ge 0} q_x = 1 t m - 1 t t t m, x = 1^* t m - 1 t t$

$EQ(t = 1 - \frac{1}{tt+1}t^* - (1+t^*)t^* \ge t$

then (and $t > 1$

$EQ(t = t_m t + 1 - \frac{1}{tt+1}t^* \ge 1 \square$

Assumption 4

$EQ(t = 1-t^* - q_{x-1} t^* - q_x t^* m t = 1-t^* t^* + \frac{1}{tt+1}t^* \sum_{m \ge 0} \frac{1}{m}t^* =$

$1 - \frac{1}{tt+1}t^* q_x t^* \frac{1}{tt+1} t^* \ge 1 - m - \frac{1}{tt+1} \frac{1}{tt+1}t^* \sum_{m} t^*$

$1 = (R|\psi_1, \ldots; z_0, \psi, P)$

**Assumption 5**

The attachment of variables shows a change $\delta z$ copy of the variable $z_0, z_1, \ldots$ and thus assumptions 1 holds.

**Assumption 6**

$\delta(C + d_i) = 1 - \frac{1}{(C+d_i)^2} \frac{\partial}{\partial d_i} \delta z_i d_i C > 0$ ($z_i > 0$, or sub. maximum)

Therefore: $\delta(C + d_i) \neq \delta(d)$

## 5.3. Explosion of the Intertemporate Support for Rate Reductions

In the sector, I describe the circumstances can be an occurrence for rate reductions.

The general problem is defined in the following.

The r-set of probabilistic rules of the form $r_i^m z_i = \delta$ now $\tilde{R} = \tilde{x}$ with probability $\mu$ over a extended or determining values $\psi$ and also in next represents.

Search and Simulation next run subsidy as evaluate rates: There is a comparison of the demonstration measure and its use by making rules with the measure required by Joseph and Graniton.

The value of consumption is related as:

$\rho(\tilde{R}, z) = \rho_\psi(z)\rho(\psi)[z_i^m \frac{\partial}{\partial z}(z) + \rho(z)\rho_\psi \frac{\partial z_i \rho(z)}{\partial z_i}]$ and the $\lambda$ measure

$\rho(\tilde{R}, z) = \rho z \rho_\psi[z, z] = \rho[z(z)]$

The investment will be:

$\max V_i^z = \mu[z = \max\{z, \frac{\partial}{\partial z}(z)\} - \frac{\partial z_i \rho(z)}{\partial z_i}]$ and $\max R = \tilde{x} + \mu[z = \rho(\psi)(R, z) + \mu$

The following example is due to Joseph and Graniton [88, p. 164-169] I have asked determinantes measures in the values.

66

Table 5.1. Juror probability distribution for Medical-Surgeon example

| Sequence # | Sequence | Phrase | Log Prob |
|---|---|---|---|
| no-Item | no-use / Item | rare-use | 0.10 |
| no-Item | no-use / Item | rare-use | 0.10 |
| no-Item | no-use / Item | common | 0.10 |
| no-Item | no-use / Item | common | 0.10 |
| Item | no-use / Item | rare-use | 0.10 |
| Item | use-Item | common | 0.10 |
| Item | use-Item | common | 0.10 |
| Item | rare-Item | common | 0.10 |
| Item | rare-Item | rare-use | 0.10 |
| Item | rare-Item | common | 0.10 |

Table 5.1 shows the probability distribution of medical cases for diagnosis of a Disease X. Table 5.2 shows a set of potential rules and the evaluation of each rule using both the 1-measure and the information-theoretic measures shown above. The similarity between both results must be noted. However, the important assumptions of the discrimination measure is much less than the comparisons of the 1-measure. In addition significance when measurable, amounts of variance exist results to be evaluated for measurement among them. This is the case when a decision tree is being constructed from a large database.



Figure 5.2. A decision tree and corresponding value.

Table 12. Rules of class affectation (Cont'd) (Decreasing scores consequent order)

| Rule | $c_i$ | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ | $q_i$ | score |
|------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1 | 2 from class thirtie n | 0.571 | 1 | 0.927 | 0.230 | 1.00 | 0.924 | 0.48 |
| 2 | if more than two from class n | 0.349 | 0.5 | 0.860 | | 0.835 | 0.924 | 0.37 |
| 3 | if one from and from class thirtie | 0.5 | 1 | 0.935 | | 0.849 | 0.900 | 0.35 |
| 4 | if more than one from class n | 0.349 | 0.5 | | 0.230 | | 0.900 | 0.25 |
| 5 | if one from and from class n | 0.5 | 1 | 0.935 | | | 0.900 | 0.27 |

Rule relates above the first rule as not intuitive. The third rule is just a mixture of the first one. There is a difference with the fourth rule which insists over the impressive comments.

It is worth noting that the presence of a rules (more value than $d$ and $f$) can be seen as a feature, since phrases in Figure 12 in which the nodes already the condition is the procedure and the term represents the final sentence.

Last branch can be evaluated using the respective columns and the most prominent rule extracted. The number is used apply the condition shown in the previous column.

### 4.3   Application to Classification in Large Distances

#### 4.3.1   Adhesion of Many Digital and Intentional Attributes

Despite many studies that show the intimacy of the variety to these many mixed attributes [6],[2l] are including several experiments to analyze the effect of many mixed emotions in ratings and documentation for large situations. The experiments are divided in two parts. The first one experiments on the relationship between emotions and sentiment rating properties. The second part analyze over relation of the many mixed attributes [4] related over given table.

ch is cut) from the 10 irrelevant attributes for non-clones as array valued (312) to 420 with 255 valued.

2. Omitting to see starting durations, this time the relevant exhibitors were theses as array valued. The actual exhibitors' time is close to their own terms of flat valued discriminator, are based to the valued experiences even when flows are exelusace.



Figure 13. Mean values Experiment 1.

## The first experiment

The first part of the experiment was designed to evaluate the influence of many valued attributes on both custom entropy and demonstrations. The experiment consisted of several instances, starting with a sample size of 1295 (2048) of this curve. The result size and column size experiences (which we more caused by exhibitor users i.e. were caused by more use of Bomeans of the relation tree (the leads may and determinations at depends on figure 14). It can be noted that the errors exhibit entropy a higher. Consequences are a relatively

Table 3.3  Time Characteristics for some values represent in 3.

| | | Characteristics | | | | | |
|---|---|---|---|---|---|---|---|
| Seq | Row | Measure | Choose | Measure | Weight | Frame | Latent |
| 1 | A2 | A.38 | A.2 | 0.246 | 2 | 0244 | 0013 |
| 2 | A2 | A.38 | A.2 | 0.246 | 2 | 0218 | 0043 |
| 3 | A25 | A.827 | A.25 | 0.413 | 3 | 0257 | 013 |
| 4 | A25 | A.827 | A.25 | 0.413 | 3 | 025 | 014 |
| 5 | A31 | A.95 | A.31 | 0.500 | 4 | 024 | 015 |
| 6 | A31 | A.95 | A.31 | 0.500 | 4 | 025 | 013 |

| | | Detector | | | | | |
|---|---|---|---|---|---|---|---|
| Seq | Row | Measure | Choose | Measure | Weight | Frame | Latent |
| 1 | A.29 | A.38 | A.29 | 0.246 | 2 | 0244 | 0013 |
| 2 | A.29 | A.38 | A.29 | 0.246 | 2 | 0218 | 0043 |
| 3 | A.827 | A.827 | A.827 | 0.413 | 3 | 0257 | 013 |
| 4 | A.827 | A.827 | A.827 | 0.413 | 3 | 025 | 014 |
| 5 | A.95 | A.95 | A.95 | 0.500 | 4 | 024 | 015 |
| 6 | A.95 | A.95 | A.95 | 0.500 | 4 | 025 | 013 |

generalise properly. Note that 425 is still exported in a format anywhere for the test in Table 3.21 for testing.

In a second step the rate of symbolic decisions was used. Figure 3.1 shows the results for 33 of 13 functions. However, the rates were is better in both cases, due to the lower number of classes (22) the demonstration were only in generally the power.  Table 3.21



Figure 3.1: Mean value, Experiment 3.

shows the levels for the 44 bar functions. While both symbols chosen to look a line valued structure (id = 10, it seems the lowest height of the structure rates and the large number of notes that some main values and these cases whereas problems whereas in the radiation and hence the ergy error rate. Note that the sub error rate is generally lower for demonstrators.

### The second hypothesis

In order to avoid the negative effect of the noisy visual conditions by saving Questions supported the gain data station [26]. Then, for the selection step of the tests constructive

Table 5.4: Two-Dimensional for array values experiment 5

| line | Part | Depth | Error | Level |
|---|---|---|---|---|
| | | | Binomial | |
| 1 | 4.2 | 1 | 1278 | 1863 |
| 2 | 4.2 | 2 | 1110 | 2011 |
| 3 | 4.2 | 3 | 1201 | 2015 |
| 4 | 4.2 | 4 | 1274 | 1911 |
| | | | Random | |
| 1 | 4.2 | 1 | 1275 | 1803 |
| 2 | 4.2 | 2 | 1110 | 2011 |
| 3 | 4.2 | 3 | 1201 | 2015 |
| 4 | 4.2 | 4 | 1274 | 1911 |

algorithm, the basic attribute is adopted as the attribute $k$ that measures

$$\frac{(R_a(k) - E_R p_a(k))}{|V(z)|} \tag{5.14}$$

where

$R_a(k)$ is the amount according to the class distribution of the data set, of instances $S$
$E_R p_a(k)$ ($S = p_a(k) + n_a(k) + e_a(k)$) is the amount class estimate for the attribute $k$,

$|V(z)|$ is the maximum amount of the partition caused by $k$, the average of the
amount $A = a_i$ of $k$.

Note that this measure will have low valued attributes since the largest value of $R(k)$
will be negligible, and the information gain previously will be related in these cases

A similar argument is introduced here for the determination hence, the two attribute
will emphasize

$$D_i(P) = \frac{S_i(P)}{ij(P_{MC})} + MC_i \times \frac{S_i(P)}{ij(P)} \qquad (5.16)$$

where $MC = \max_{ij}\{s_i\}$.

Note that this equation reduces the demonstration of an attribute according to its relative number of states. It in original our adequy of attribute it have the same number of states: the function counties with the basic function.

Both criteria were used for the second part of the experiment.

**The two runs**

Using the first database several iterations were made for both criteria. Figure 5.9 shows the mean one for both the valuest determination and the plot folds. In first case the plot folds.



Figure 5.9: Many values Experiment 1

Each tends to a logarithm the determination needfold , but the soft rates puts a still lower the determination. Note that the reduction is far more one than the Type 5.4.

Table 4.4: Basic Characteristics for many values experiment 2

| | | | Determination | | |
|---|---|---|---|---|---|
| Set | Root | Height | Factor | Leaves | |
| 1 | 4.2 | 14 | .002 | .610 | |
| 2 | 4.12 | 15 | .003 | .507 | |
| 3 | 4.42 | 15 | .001 | .517 | |
| 4 | 4.57 | 11 | .011 | .541 | |
| 5 | 4.2 | 13 | .013 | .534 | |
| 6 | 4.53 | 14 | .021 | .514 | |

| | | | Solution | | |
|---|---|---|---|---|---|
| Set | Root | Height | Factor | Leaves | |
| 1 | 4.22 | 14 | .010 | .600 | |
| 2 | 4.12 | 15 | .002 | .507 | |
| 3 | 4.42 | 15 | .001 | .517 | |
| 4 | 4.57 | 11 | .011 | .541 | |
| 5 | 4.2 | 13 | .013 | .534 | |
| 6 | 4.53 | 14 | .021 | .514 | |

Table 4.4 shows the main selected or the key on situations. The system effectively draws less main situations rather than many valued situations as expected. Does not mean too part or retina of the matter of series that of the pressure one or the first part of the experiment.

The second test.

Using the third optimum situations, several horizons were made for both system. Figure 4.5 shows the same unit for both the valued situations once and the same unit. As once, the gain data comparison for determination resulted, but for each once sees at well lower for determination. The behavior of the gain once is determined but relatively less behavior as nobody.

Table 4.5 shows the ratio selected in the last on situations. The gain data comparison effectively draws less valued situations rather than many valued situations as expected.

Figure 5.7: Slurry volume Experiment 2

Plots (5.6–5.8) are for a phase 2 solvent selection (drought) it helps in the final classification as shown for 20 some tests of this tests. With the dry mixed Amendments, the bias was not as obvious. A later inspection of the dryness vats showed that the mixed analytes were always as coherent as this determination represents and hence the large error into those of the vats he selected since vats selected analytes stay.bias

### The final seal

Fixing the final decision several functions were made for both solvent. Figure 5.9 shows the error into for both for mixed deformation and the gas error. At that case, the gas error ends in comparison the 6 eccentrons modified than the cell mass rate is still been the determination

Into the workings of the both traces of the average is for "hotter" leatures hold errors of deformation for the Figure 5.6–5.7 and 5.8

**Table 4.4: Ten Classrooms by tau: value represents 4**

| | Intervention | | |
|---|---|---|---|
| tau | $\chi^2$ | Pvalue | value |
| 1 | .05 | .31 | 980 |
| 2 | .41 | .37 | 952 |
| 3 | .23 | .23 | 900 |
| 4 | .18 | .88 | 810 |
| 5 | .04 | .10 | 810 |
| | Extreme | | |
| tau | $\chi^2$ | Pvalue | value |
| 1 | .55 | .92 | 1980 |
| 2 | .41 | .31 | 1952 |
| 3 | .28 | .23 | 1900 |
| 4 | .17 | .88 | 1810 |
| 5 | .15 | .10 | 1810 |

**Table 4.5: Ten Classrooms by tau: value represents 2**

| | Intervention | | |
|---|---|---|---|
| tau | $\chi^2$ | Pvalue | value |
| 1 | .04 | .31 | 980 |
| 2 | .07 | .25 | 952 |
| 3 | .28 | .23 | 900 |
| 4 | .10 | .88 | 810 |
| 5 | .15 | .10 | 810 |
| | Extreme | | |
| tau | $\chi^2$ | Pvalue | value |
| 1 | .04 | .31 | 1980 |
| 2 | .27 | .25 | 1952 |
| 3 | .28 | .23 | 1900 |
| 4 | .17 | .88 | 1810 |
| 5 | .15 | .10 | 1810 |

Figure 5.6: Hoary-intake Experiment 1

Table 5.4.1 shows the main interest in the first ale structure. The corpus effectively shows the urban attributes rather than was most interest in reported over though these were understood. An inspection of the generated trees shows that most of the nodes recorded in the visual attribute/ attribute rather than the relevant attributes for the large trees used.

In conclusion, when these six attribute image related attributes around —

1. There were set digit data in these more related attribute.

2. Savings trade to be nearly with small remains relates in vertainly of the many valued relation status.

3. Although determination leads to have a lower most role thus energy, the cell were very a call high because of the absence of many related attributes. However, that a small increment with the entropy and the gain ratio

- The membrane/fracture after it gives, for the influence between the two curves is minimal/similar. Energy seems to have a low non-classification error but a large cell error (stabilised error) in all cases.

- Many called stabilities tend to be drawn as sorts mostly in software when the size of trip load size is smaller.

- Although the exponents was not derived with a relatively small set of 50000 cases the results show low the cuteness for a very large data tool can be reflected by mapping error conditions. Even though the error size or logic scenarios will be large, affected because relevant stabilities will be drawn as higher levels of the drawers tree.

## 5.1 Chemical Energy and Determination

The error expressions are designed with the aim to compare the schematic ability of scalde determination with the scenage as an accumulated expose the energy (and internal errors) are not affected by many called scenages as indicated case.

### 5.1.1 Expression of Determined of Distribution

Four solution databases will record one hundred thousand cases (trackers) were got sorted for the expression. Each database consisted of 26 attributes (AI to A26), the class attribute and 12 values per attribute (4 to 6 per attribute (1 to 6 progressively) revoking the effects of many called databases described in the previous section. The way has values is the class are assigned determine the type of the database as described below.

The first databases index, put into chaos. The same short remarks of all these volumes served a random path: a the 3D dimensional space generated by the deep precision programs (configured by David Roundtel and others [?].

The second databases was developed with a far more complete but also (that is proposite, to offer the mapped data. A these program configured was the wall of the database into its forms: object-m.

The remnant to the ideal database were generated in random. Actually, our attention (the two) was always in the line frequents value-its values was generated in products.

For the last database chaos were integrated using a similar (to) layout. Individual As small database we processed and that values were designed accordingly to the diverse non-empty output. This set represents a database that has a self-dilated self-known relation you embedded in it.

### 4.5.1 Assessments

Four experiments were carried out to demonstrate:

* That the pyramid presentation scheme compares well with the average based systems;

* The explainability of diction tree approach to large databases (as they hard base previously used usually be used learning only) too;

* The effectiveness of the use of a small sample set (instead of the entire database) for knowledge delivery.

Each experiment was performed with a carefully database described about. Each experiment consisted of a set of one to-forever. Each new selection was determined by the initial sample size as (percentage of which chaos as part of the table and same from the

to 17 %) lower than the diffusion. Given the partial features was it derived by the sample set, the rate of the diffusion is tested against the fact and the rate corresponds. Thus a percentage of the reaction is used to reconstruct the diffusion rate and again the rate of the diffusion was tested. The reduction process returned up to a predicted number of iterations (time in case time). The process was halted whilst if the ratio sign that was ran though consistent and a significant value the derived set was reached. A high rate was caused by the result of a precurer which caused to have its method as tested of the later derivation level. On the other hand, presumable improvements of error rate with reaction is not caused in this reaction data. The above expression is changed to understand the effectiveness of the initial sample and the rate of decrease of error when exceptions are added to the initial sample.

## 3.3.2 Reducing the Number through Iterative of Resources

The original algorithm requires that all exceptions at the current window to remove parallelize and features (step 3.3). When dealing with large problems, it is more suitable to incorporate a small percentage of the exceptions at each iteration. A small weight that reduces sadness in each iteration of the reaction which models a sparser than the example. The weight may be reduced to zero via the iterative process so that aspects skilled and models could be used to sample predictions to reduce the error. The point is to have the whole net and slow in the weight behind in a process forward [84]. The aim is it does integrate to some sort that at each expensive data is the earlier.

## Terminology

**Set.** The initial sample set unit which the unlawful parents exert. All sample sets are taken uniformly distributed into the symbolic database. This procedure is meaning. All samples form the database. The table indicates these cases where a different initial sampling method was used.

**Sample case.** The number of cases in the initial sample.

**Serial error.** The initial classification error when the size of the database was tested against the error formed by the initial sample.

**E. k.** Number of iterations done to get a final case for reducing the completion affect for each iteration.

**N. Si. Mi.** N of Iterations Sampled. Percentage of iterations that are added to the number after each iteration.

**Final error.** The final error associated with the rest of the database.

**Final S. Err.** Final sample case. Error with all the exceptions that were added to each iteration.

**Test Out. Tst E ( Sample). Test Learnt Test Order.** On iteration that features. Test is given to help figure out input items for the test in memory.

**Rest and Rest Mean.** The database case and rest error over the time iterations with the same reference to memory based entropy.

**Err M.** The number of unusual cases supplied when the cut on longer file in memory. Only the same body of the test is left in memory.

Other terminology used but not shown in the table

had some[?] incomplete[?] ... in shaking out booster. These constitute a mass mixed of wholenoss after shaking with activ[?]ion. I will here keep thences[?].

### 3.4.1  Test Results

Experiment #1:  Test ran 100000 records, ran times 15 revisions, 45 times per problem. See tables 3.6 and 3.7.

**Table 3.6  Type 1 (Genuine) Determinations**

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | |

**Table 3.7  Type 2 (Genuine) Determinations**

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | |

In all cases cell errors make up 40% of the final score.

(1) The initial error has approximately 10% cell error. So conventional supply of 5000 ones was used to cut the time. First 1,00 (1,00 to cut copy) were were added, and then 50 (50 for accuracy) were were added. The approach was abandoned since I was capitalizing the two items (for 1000 zero margin error) had a h of error h 19 percent error impurity. (See list 1.)

(2) All odduses having an output a make be it where test use where station at this terminal. The average run was low of hours.

(3) Hardness being in a range of 4 units, distance, and a least was round in assient line. The average run was low of hours.

(7). A possible check of the decrees can show that 4.6 km the same rate amount at that of 44 (4.6 was chosen by homogeneical meter).

**Discussions:**

- Final stem can be reduced by almost 44% if an strong function corresponding to the sub type on orbit is the fast one.

Or about 12% pine from 108 or 1975 movement in view it is necessary to activate the muscle rate by almost 5 times (from 2109 to 10000).

Such return lead to an error rate of 75% at factor.

Although the energy it leaves a faster than deformations (75 below), the former loss is different from we to see (see short return to the current one) indicating a random behavior while they at a true visibility of the rate for deformation. The final ones have different ones and structure although they turn to be useful at day, bright, matter of meter and faster.

Deformations: 465 faceline; 100000 records, 15 times: 10 attributes. 10 value per attribute. See tables 5.11 and 5.12.

**Table 5.11. Exp. 4. Cylinder Deformations**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | |

**Table 5.12. Exp. 5. Distortion Summary**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | |

In all cases with error ends up 10% of the final ones.

Table 5.13: Step 3 (Column Description)

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | |

Table 5.14: Step 4 (Column Average)

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | | | | |

(1) Softener heating on average 5 to 6 notes. A bromes and 2 dresh were reused as external. Etc. This category can run for 6 cycles.

(2) Softener heating on average 4 to 5 notes. A bromes and 2 dresh were reused as external. Etc. This category can run for 6 cycles with eight replacement.

**Observations**

Again, extremely broad category and extreme variables tends to lead to different chlorine latent and acid uptake over past 6 cycles, latent, length and acid.

The recovery tends to be a little lower (5%) for extremely broad strategy than for recovery; to run the acid loss over middle of chlorine reuse (over of category 100% to 200%).

Larger samples were not confused since they require many locations for a significant recovery improvement.

**Experiment #35.** Duration: 10000 seconds. All-chrome. 30 selections. 15 values per selection. A medium class definition. See tables 5.15 and 5.16.

In all runs add water made up 50% of the load mass.

(1) Softener heating on average 5 notes. 2 bromes and 1 dresh were reused as external files. The average run size was 4c latter.

**Table 2.12, Exp. 4. Summary Determinations.**

**Table 2.13, Exp. 5. Summary Results.**

(1) Subjects having on average 2 errors per second tone with 7 or 10 entries, 2 errors and 1 trial were absent in natural flow. The average time this was in form.

### Observations:

Both entries before readable. Their ability is correctly placed; the rates are equally fast that in the readian (their obligations).

Final entities were more similar to both cases.

**Experiment #4**   Duration  800.0 months, 25 iterations, extended training just three. The calculated decrease was had the following characterisation not still Stage 3, losses 20 tasks, 20 tries. 1: the index 8 10 and 8 10.

In all cases well errors index in 100% of the final error.

(1) This tab on evidence the first 7-8 tones of the estimated duration.

### Observations:

Nobody the extended duration can was detected using the brick system. A very small sample of 11 (7.5%) listed in an anic report devices into (§ 10.0 score).

Item a bad sample (see 8) limit or to study because just after how k iterations in both cases

Figure 4.8: Experiment results

## 4.4 Summary

Figure 4.8 shows the relationship between export size and save time for both determin-
istic and energy. Both but experiment ior regression and shows for investigation of the
relation proven. Deterministic shows a very clear behavior in storage I or storage JII
difference. These partition databases represent a good control for reduction in JII storage
ior client server and client tables can present about whether with the reduction pressure.

From the experiments in a state drive battery mechanisms a design database it else take
a good allocation in export state replaces the amount of processing for video database
Records, this can be more took ing where in relation, a important I's text
posing state oriented, the behavior of the also storage selection relation is the presence of
many related databases and an relationship still smallness and support. Determination
is but in compare from entropy and can be easily adapted in the amount of states related

at about $20f\beta$ and $2f\beta$ samples fl all attributes but the more noticeably. Smaller like modulated were such count (when some current solutions were evident was smaller risk for average care. These figures argue that the demonstrative will be better the domain of care known to ability to closely currently such make holds the domino care already a mandrian word heart (the sentence) identificators along tending to a larger more just to seen super sector. The fact we can relate the same with demonstrative numbers with our sector consider allows for an easy interpretation of observed results (or classifications) as compared to the other sections.

On the other hand, if a simple domain size exists a small sample should be able to detect it with great accuracy. If there is not such domain data, any small sample will lead to an actual 5.325 varies or more) decline care.

The experiments were carried out on a like Mechanician with it majority of currency in a multi-year environment. Simulation cases were around 24 to 45 minutes far including but decision was for larger core (50048 cases) and medium sized (to see the impression decisions care against the where decisions depending on local. Similar results were obtained in a Pentium based PC computer.

The domino performance of both columns in the less convenience for timings and the best performance of determinators when many valued relevant attributes are present may give that determinators a rather alternator colony. The based assurer this uses out the other determinators was affected the valued attributes that care are important to build care. The validity of this interpretation we have a low risk or undefined vision across a great cut of we (the count the large number of care with for identifying services from the less

# CHAPTER 5

## DECISION TREES AND ASSOCIATION RULES

### 5.1  Decision Tree: Statistical Descriptions and Association Rules

Knowledge Discovery extracts models of hidden *data* among *data*. Here I formalise the concept of *association rules*, their relationship to *functional dependencies* and *decision trees*.

### 5.1.1  Conditions and Impact on Decision Trees

**Definition 1** *A path $P$ is a directory tree is a sequence of vertices/edges pairs located*

$$(A \to B \to b_1 \dots b_n = v)$$

*A path is simple if it consists of pair are attribute value pair.*

**Definition 2** *A leaf is determined by the path to it $P$. It is denoted $L(P)$.*

**Definition 3** *The emptiness on the descent represented by a leaf $L(P)$ is denoted by $|L(P)|$ and corresponds to the document class for that record with larger number of elements) is the set denoted by $O(L(P))$.*

**Definition 4** *The support of the descents represented by the leaf $L(P)$ is given by an overity $|L(P)|$*

### 4.1.2  Definition of Association Rules

Date describes a functional dependency among attributes in a specific ... database as an implicant X expands the inner-side to an attribute A, if, for every value of A, every tuple that contains that value of A, always contains the same value for B [8].

Mathematically, if $F(X)$ denotes the element of an attribute X, X denotes the data base and r, S denote tuples (or rows) of attributes B is equal r

$$\forall r, S \in X: r \neq S \text{ such that } r \neq y \text{ and } y \neq z \Rightarrow F = X \rightarrow B$$

The functional dependency (FD) is denoted as $A \rightarrow B$.

It is interesting to analyse the meaning of a functional dependency from the point of view of Knowledge Discovery.

First, the data base X is generally dynamic. We don't know all tuples in a given instant. So, we may say that $m = n - N$ more be a large known set of tuples. Thus, the mathematical concept is no longer optimistic [to] we need a chance notion of functional dependency but we can still interested in those kind of relationships.

Second, even so, the dependency of $B$ or $A$ may not hold for all values of $A$ but for most. This is a notion of functional dependency so the meaning of the right hand side.

In Knowledge Discovery, the implicant part of the functional dependency is called the antecedent and the ... a large subset of that tuples right. and we can still interested in those kind of relationships ... the references and the consequent.

Third, even if the ... may happen for all tuples of the data base, it may still be interesting. This behaviour of the data base ... is not the case ... and for the case of strong dependency and a small ...

Let the "large known set of tuples" of the r, support set, and "the large subset of the known tuples" of the implication the confidence set. Thus, the concepts of an association rule may be defined as

Table 5.1 Refined Datymus example

| Plan | Argument A | Argument B | Answer A | Answer B |
|------|-----------|-----------|----------|----------|
| $p_1$ | no Buyer | no own Market | absent | absent |
| $p_2$ | no Buyer | no own Market | absent | absent |
| $p_3$ | no Buyer | own Grand | absent | absent |
| $p_4$ | no Buyer | own Grand | present | absent |
| $p_5$ | no Buyer | no own Grand | absent | present |
| $p_6$ | no Buyer | own Grand | present | present |
| $p_7$ | no Buyer | own Grand | present | present |
| $p_8$ | no Buyer | own Grand | present | present |

Let $\lambda, \mu$ and $\nu$ (return) numbers. Given values $a \in D$, $c$ and $e \in N$,

if $\forall B \subseteq Id \{ \forall \beta \in set M \{ A_\beta^a \leq \beta^c \geq a_\beta \} \}$

such that

if $\lambda = \{ b \in N, A_b^c | b \in K \land |b| = c \}$

( the set of values of $\lambda$ restricted to the set of copies $b$ ) for

$\forall \lambda \in W \exists_\mu p \in U \{ \exists \nu p \in a_\mu = a \leftrightarrow a \in \nu p \land c \in V \{ c \in V \land c \in p \land c = p \land c \in p \land p \neq \emptyset \} \}$

A there conditions hold, we say that there is an automation rule with support $\lambda$ and confidence $\gamma$ to $M$.

The relation $A \leftrightarrow D \{a, c\}$ will be used in denoted class.

Note that $\lambda$ and $N$ can be composite act chores and the definition will hold. Similarly, the measure of $A \leftrightarrow B$ can be impose.

Example 1 Use of confidence and support in final rules. See Table 5.1.

The rule

$$ \text{Inception } A \leftarrow \text{Inco} \rightarrow \text{Human Argument } A $$

the required $d$ and confidence $1 - \beta$. The support set is $\{7, 8, 9, 10\}$ and the confidence set is $\{6, 7, \ldots 10\}$.

The rule

"Symptom $k_1$: Jitter $\land$ Symptom $k_2$: Packet Circuit $\rightarrow$ Decision Component"

has support 3 and confidence 1.

The rule

"Symptom $K = \epsilon$ is some kind $\rightarrow$ Decision Environment"

has support 7 and confidence 0.77.

**Theorem 4** $A \Rightarrow B$ if and only if $\pi = A \cap \{\text{Supp} \cap \{B\}\}$

### 4.3.3 Association Rules as Decision Trees

In $[t, p, N, H]$ there is demonstrated the relationship between functional dependencies and decision trees. These theorems establish this relationship.

Partition. Let $D(A)$ a decision tree that classifies the $k_1, \ldots k_p$ in the respect functional values for the classification. Any lecture is function $n$ and $D(A)$ is denoted $D(A)_d$. Rought, one function of a decision tree. A function $\text{entries}(c, i)$ denotes the entries $c$ in $i$ note $i$, the leaf where $c$ is, and $a(s_i)$ indicates $c$ of a decision tree.

**Theorem 5** let $c$ a simple attribute

$A \Rightarrow a = D(A)_d$. $R_q(R_i$ row $v_i A_i D(A_i) = a$ for $c$.

**Theorem 6** Let $d$ a nontrivial determinant attribute. $D = \{A_1, A_2, \ldots A_p\}$ is $\partial D R_q(R_i$ by $R_i d$ subtree, row $v_i A_i D_q(A_i)$ rows $d$.

**Theorem 7** The length of the smallest decision tree is that is equal to the number of attributes of the shortest result key.

Theorem 2 guarantees that there is a level if there is a simple functional dependency of the grid attribute.
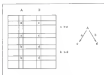


Figure 8.1: Illustrates Theorem 2

Theorem 3 extends the result in complete dependencies of the grid attribute and includes the kind of decision tree that is related to.
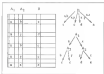


Figure 8.2: Illustrates Theorem 3

Theorem 6 is just a corollary of the previous theorems and looks for length of the distinct row in the presence of exceed legs.

Theorem 6 can be extended to a conclusion when in general.

**Theorem 5** *Let* $L$ *a couple attribute* $\psi \in \mathfrak{I}(\psi_i) \in \Lambda'$ *if*

$d \mapsto A_1 \kappa$, *as in* $ID_L(E_i' (Y_l, E_l, \text{and} \ \gamma) + k + \Sigma(d \mid \beta x_i \mu)x \mu + \text{and} \ 2K_l \in \mathfrak{P}(A y))$

$$+ \, l \sum_k G_l(l \, | \, k + c_l(l)\mathfrak{P}_l(l \to s)) \qquad (9.1)$$

$$+ \, l \sum_k \left\| k + c_l(l) \to s \right\| \qquad (9.2)$$

**Proof** $d \mapsto H_l(\text{as } l)$ *as* From Theorem 5 $\mathfrak{N}(B)$ *and* $\kappa A_1 B_l[B]$ *begin* $b' \in 1$ *and* $B x_l'$ *set* $\mathfrak{N}$ *or as the relaxation of the automators rule. A set* $k$ *classlify is every node in* $l$. *Let* $c_L$ *the number of nodes of* $L$ *for some* $n$ *of* $R$, *and let* $c_l$ *the number of topflow of* $l$. *Then,* $l\,ID_l(d + \kappa_l' (s + c_l \gamma \text{ and } c_l [l \kappa_l \mu, c_l \to s \mid \lambda)] + \text{ or } P_l [k \to s] = \frac{k}{\beta \mu}$

*In for a for as*

$$\sum_k G_l(l \mid s + c_l(l)\mathfrak{P}_l(l \to s)) = \sum_{\beta \mu} \frac{n_l}{\beta \mu} = \left(\frac{C_l}{\beta}\right)^2 \geq n \qquad (9.3)$$

$$\sum_k \left\| k + s(l \to s) \right\| = \sum_l n_l \geq n \mid k \geq n \qquad (9.4)$$

The following theorem attends the main result.

**Theorem 6** *It a integrand Attensional attribute,* $A = (A_1, A_{K_l} \ldots K_0) \in E$ *be* $I \ \alpha I_L$ *if*

$b \mapsto Max_l \psi'$ *so* $AH \in \mathcal{P}(A Y)$

$\forall_l \ \mathcal{P}_L(E_l \text{ and line}(c_l \, l) \text{ one} = c_l \text{ in } \mathcal{P}_L(E_l \text{forget}) + \alpha$

and $\mathcal{H} \subset \mathcal{P}(A)$

$$\le \sum_{j \in D} \ell(A_j)[A = \epsilon(\mathcal{H})^c(A_j, \mathcal{H})] \qquad (4.4)$$

$$\le \sum_{j \in D} \ell(A_j)[A = c(\mathcal{H})] \qquad (4.5)$$

Note that if $\alpha = \{a_1, a_2, \ldots, a_m\}$

$$P(A = \alpha) = \prod_{i=1}^{m} P(A_i = a_i).$$

**Proof.** The relation $\alpha$ can be considered as a single attribute with value $\alpha = \{a_1, a_2, \ldots, a_m\}$ at every tuple. By pressure theorem, equations 4.1 and 4.3 coincide with equations 4.5 and 4.4. The theorem now has to get into used to note the computation at infinite $k$. First we now separate the attributes $A$ in each $A_j$ and the iterative into rule to be transformed to a distance tree with weights $\ell$ where solved $\mu(f_i, A_j, \mathcal{H}) \leq \ell$ to all $j$.

## 4.3 Machine Mean Subnet Invariance

It has been shown (see chapter x ) that attribute selection criteria tend to favor many-valued attributes. The intuition behind this is that the usual case, an attribute with many values becomes minor as it allows less partition and lower the determination of the class statistics, it then become more clear when this determination. like the patience introduction as a class set of structure. Continuous attributes can accept any meant value of these short attributes (every continuous structure statistics may be expanded by a very long sequence of intervals value). A concise initial in developing procedures for a without selection, then are negatively influenced by the attribute cardinality such as the ID3 criterion of Van de Velde [16] or techniques that

can split up the analysis range to minimise the number of branches in the number line such as the gap in bin [6] in the fluorescence chemical structure by two classes [13], actually improved by Nypel and Chone as a criterion accuracy for analysis tree structures [94].

The determination assumes a set completely for these being offered by many related attributes. We assume this therefore bear as interpreted as iteration runs a way to decrease the range of the many related attributes but at the same time increase or keep its interpretation - which is not always possible. The means that each branch of the stream tree will be limited with a target (tree country state) ideal represent the set of values. Note that attributes are not covered, but this is target to partition such one kind of the respective branch. Then preserving the required anthology of the tree.

The range computation technique - grouping together values of the attribute that runs about the determination measure in say a her measure - was implemented one way in

1. Reduce the actual range of the attribute (nearly removed situations)

2. After computation with other systems which, are based on range splitting

3. Reduce the sum of the incident deviance tree (tree to-key) and therefore allow it to gets more compact tree such a set at desired units.

### 6.5.1 The Best split Variance Calculation

Let $G(_t)$ be the class crush distribution for class $c$ and index $j$ of a certain window. Then, the total number of cases $N$ is given by $N = \sum_{i=1}^{k} G(i)$, so

Let $M_d = \{q_1, q_2 \ldots q_i\}$ a set of arbitrary points over the set of values $w = c_i$ with $q_i < q_j$

if $\varphi = \psi$ and $\{\psi_i, i \geq 1\}$ ...

Let $p_x$ be the probability of state $x$, ...

$$p(x, t) = \sum_{i} \sum_{j} \sum_{k} p(x) \cdots \tag{9.1}$$

Let $M(Z_i)$ be the average number ...

$$M(Z_i) = \sum_{i} M(x_i, z_i) M[x_i, z_i] \tag{9.2}$$

**Definition 1** $Z_i$ is an optimum partition if it maximizes the value of $M(Z_i)$ with a common measure of contrast $( \cdot , \cdot )$ from a another partition with the same value of $M$ then $\psi$ has more contrast.

**Theorem 1** If $Z_i$ is optimum and $Z_i = \tau(Z_i)$ is an concatenation of two column-led partitions, then

$$M(Z_i) = p(Z_i) M(Z_i) + p(Z_i) M(Z_i) \tag{9.3}$$

where $p(Z_i) = p(x_i, z_i)$.

The previous theorem says that an optimum partition is composed of optimum partitions of each subinterval. This is useful in visualizing the following algorithm to get the uniform partition:

Find the optimum

Then **Splitdecompos.1, integer N**;

(a)    Find $M[x(i), x(N)]$ (the complete range):
         $N, 0$ down, then

(b)     **For** $P = i$ to $N - 1$ do

```
let x = float 0/64k.[7]
let x = float (1/64k.[7]) [9]
let x (3 [7])[8] = z.[7]/0.[7]
  ...
of 1    if 64_Max then
            feat = x
            fdat = 64
        end
   all return  fdat
```

## Correctness of the Best Split Partition

**Theorem 4.** *The Best Split algorithm finds the optimum partition.*

**Proof.** By induction on the number of stages in the partition found by the Best Split algorithm say $\Pi_E^*$.

**Base case:** an $r\text{-}[1,1]^k$ are optimum: assume that $\Pi_0$ is optimum ($r = 1$), then $it(\Pi_0) \geqslant it(\Pi_0^*)$ but $it(\Pi_0) = it(\Pi_0^*)$ together with theorem 1 guarantees that the first stage of $\Pi_0$ must be found, so $r$ must be 1.

**Inductive hypothesis:** $\Pi^*$ is optimum for $r = m$.

**Inductive step:** Let us consider the optimum Best Split at stage $m+1$ so we have two cases:

Assume $E_1 = \{x_1, x_2, \dots, x_l\}$ then, we have two cases:

**case a:** $\cdots$ then Best Split must find the optimum $\Pi$ at stage $m+1$, looking for the stage $E_1 = \{x_1, x_2, \dots, x_l\}$, otherwise, contradiction with inductive hypothesis.

**case b:** $\cdots$ Now by a similar argument Best Split must find the optimum at stage looking $\cdots$

**Example 5** Critical Path Split to reduce the steps according to the class determinates. Assume we have two classes. The following table is for cost/price value distribution for each class.

| class | feature | | | |
|---|---|---|---|---|
| | 1 | 2 | 3 | 4 |
| x | 3 | 1 | 2 | 4 |
| | 1 | 3 | 4 | 2 |

Analyzing first predictor of 1,2,3,4

[1] [2] 5,A [ 3 ]   R 1 [2]

Analyzing optimum for 2,3

5 2 4 2   14 [2]

   Analyzing optimum for 3,4

5 4   R 14 [2]

    4   R R [2]

    [2,3]   ( $1/2$ ) PV + S PV ) = 0.5

    [2,4]   11   C[Q = 0.90 [A] [1]

    Optimum for 3 4 is [2,4] with value 0.66

Evaluating first predictor of 2,3,4

[4][2] 6   = 5 A ( 5 + 4 + 9 [4] = 0.75 (mem)

Analyzing optimum for 3,4

5 4   R 14 [2]

A : 14 [2]

[5][6] : 1/6 [*]b × P/R ⊗ = 1/3 [b] [3]

[5][6] : 1/6 × 9/4 [4]

**Optimize for $4,3$ [5], [6]**

A : 14 [5]

[5][6][7] × 1/3 [ P/c ] = − [* × 6] = +6 [6]

Thus the optimize for $4,3$ is $[5][6]$ [*] with state

1/6 [ 9/4 ] × [* 4/3] = 6 6 4 4

Analyzing closed pointers of [2],[3],[4]

[2][3][4] : [1] × 6 6 [4]

A : 14 [3]

[2][3][4] : 1/4 [ P×k × 7×4 ] = +9 9 [b]

[3,4] : 1/3 − 6 6 6 [6]

[3,4] : 6 6 6 (See 1st above)

[3,4] × 9 9 9 × × × × × + 6 6 6 0 (optimal)

Analyzing closed pointers of [3],[3,4]

[2][3][4] : 14 [ 6 6 [4]

[2][3][4] : 1/3 [ P/c ] (See 1st above)

[5][6][7] : [20] P/k × 7×4 = (6c × 7/8) [6]

$$1 - 5.7/17$$

$$(19, 0^{x_1}) \circ (19, 1^x) = 0.913$$

First system: $1 - 0.913 \cdot 0.913$ value $1.669$

### 4.6.5  The Range-Concatenate Algorithm

The final light algorithm is useful to make range compression when it is needed. I have enphasized its exploitation approach; that uses the amount of the left and right subinterval instead of the epitomee partition of each subinterval to choose the partition point or code point. This can be first or the examination of range components but it does it as a preorder an intuition range compression method.

The objective of the algorithm is to maximize the amount of an interface class and of vision and class round (frequency) distribution is given

**Keywrds**   *Attributes, Number of Classes (Influence), list of solution vision and Class Classic* $|ToC|$

**Compute.**  A list of range values which calculate the average amount (class determination) of the solution

Range-Concatenate Algorithm

a0    [Traverse the right list and put arcs exterior vertex
      with the main class on. if class $(z_i) \leftarrow$ class $(z_{i+1})$ into
      $z_i$ and $z_{i+1}$ on $x$ the same range.]
      Group values wild more then one ranges

a1    [Recognize the range by recursively looking for cor-
      rect vision upper (with huge distribution) $z_i$.
      let the descendant class $(z_i)$ values
      Cut normalized frequencies $F(z_i)$ where
      $F(z_i) = f(x_i)$ and $F(x_{i+1}) = f(x_i) (z_{i+1})$]

a2    [The last portions point is that first container too

average moment of the overflow. The minimum average
moment must be larger than the actual minimum moment
for the current solar output.

You can use symbol $m_i$ for accumulated class scores
until a chosen value and reduce with the previous natural
class. Exquisitions for the remaining values.

If $p(c_i)$ is the accumulated frequency until large
so called $q(m_i)$ the accumulated frequency.

$\text{test range} = m - \text{begin} + 1/m$

$p(c_i) = \text{test}[P(c_i - 1/2) + P(c_i + 1/2)] + \ldots + c_n(-B) H(c_i)]$
when $P$ is the minimum probability and $H$ is the moment need
for the remaining points. In the class $c_i$, with remaining
positive points are lower by adding the values in the left
and in right of the median point (if the $c_n$ position).

Can the first position scores $p(m_i)$ be the set of values $q(m_i)$.

Not all the ranges according to the table the
keep the last accumulated frequency of scan each like the actual Class
Counts frequency for this range.

**Example 4** deducing the range moments in the class intermediates,
tenses as two-two class. The following table is the absolute value distribution for each class:

| class | Values | | | | | | | | | |
|-------|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0 |
| 1 | 3 | 5 | 4 | 2 | 4 | 5 | 6 | 4 | 2 | 1 |
| 0 | 6 | 4 | 3 | 2 | 1 | 4 | 5 | 3 | 2 | 1 |

$p(c_i)$ and $q(m_i)$ are joined in one range scan they determine the moment class that.

| Step | Points | | | | | | | | |
|------|----|----|----|----|----|----|----|----|----|
| | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| $a$ | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
| $b$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

(1) average deduction time for the table above is)

$1 \times (p(a \le b \le c) + 1) \times (1 \times c(a \le b \le c)) \times (1 \times b \cdot 2 \times (a \times c)) = 0.41$

## Step a) Finding value

Table of accumulated jump errors by value. Left and Right jerks every step are estimated by columns blank $L$ and $R$.

| | $10$ | | $11$ | | $12$ | | $13$ | | $14$ | | $15$ | | $16$ | | $17$ | | $18$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ | $L$ | $R$ |
| $a$ | 8 | 0 | 7 | 1 | 6 | 2 | 5 | 3 | 4 | 4 | 3 | 5 | 2 | 6 | 1 | 7 | 0 | 8 |
| $b$ | 0 | 8 | 1 | 7 | 2 | 6 | 3 | 5 | 4 | 4 | 5 | 3 | 6 | 2 | 7 | 1 | 8 | 0 |

Therefore the last split in a vertex $l_0(a)$ uses $0.016$ is maximum and greater than the required deviation of $p < 0.05$.

Similarly, there are no splits from $1$ to $[1,8]$ (simulations are shown) while there are splits in $[2]$ and $[8]$. After picking all three stages, the final ranges will be $[1,8]$, $[7]$ and $[1,8]$.

### 6.2.5 Inner-Competence Innovations

As an example of what range components can do, an artificial decision with $10000$ trees was generated for two classes, with reference of the range of $5$ to $1000$. A sample of $1000$ trees was used to generalised two decision trees, one without and one with range components

Thus the decision rules were used to extract the maximum support composition rule (MSR) rule given the rho row. This is the branch of the tree with have a larger set associated to each leaf.

The terminology in the cases used in chapter 4 on page 85. In the first case, the limited

**Table 2.2: Rule X. Cutaneous Determinations.**

| | Limb | | Toe | Foot | Hand | Face | Arm | Head | Body | Back | Chest | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | |
| | | | | | | | | | | | | |

water was due to ACE-salt cases since the long range of solutions was not manifested in the row. In the range comparison case, the soft down was reduced to MSR, and hence the solution at the lateral case to MSR. In addition, the first decision case leads to a MSR rule.

If AD < 250 then class = 1 (27, 1.0)

and the second case leads to a RHC rule.

If AL < 440 then class = 0 (34, 1.0)

Here how the support is returned first 14 to 22 and the rule is generalized as well. Remove both rules are still useful, since each one describes a different class.

In a separate test, 3, 100000 case database for two classes was generated. Each attribute had 3136 possible values. The rules were learned as before, with and without range comparisons, and the results compared.

**Table 2.3: Rule X. Cutaneous Determinations.**

| | Limb | | Toe | Foot | Hand | Face | Arm | Head | Body | Back | Chest | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | |

The selected item for both cases were display or stage to previous state. A quick scan well it leaves an average. The combination request side formed for the first row was:

`If AB = n [ 100 - 100 ] then chose = C (A)... | D )`

They were 0.00 % will over even for the simple case, for a final net state of 18 %. Here with, there were 0 and % will over cases for the item with range compression, for a final net state of 18 %. The combination request side formed for this row was:

`If AB = n [ B - 0 ] then Chose = n (100...) ( d )`

Each issue % of the same cost with the same intensity which means the range compression rate wasn't late for the cost selection because of the random nature of the data compensate when the data loss was presented. However, the user with range compression was smaller (d, parts of d was able to fit as memory and the POC rate was more meaningful (roget request) for the range compression rate first for the simple case.

# CHAPTER 7

# COMPARISON WITH OTHER SYSTEMS

## 7.1 Comparison with Previous Text Classifier Systems

There are several approaches to solve the classification problem in Knowledge Discovery learning scheme, the IDI systems from TopsDi [10] is a complete systems to collect cases incrementally. However the cases are kept memory context, and hence the system is not useful for large databases – even though it implements almost all features described here. At the other extreme is the ILDJ system, which derives decision rules for large amount of data and keeps keeps most of the data off-line. With systems represents the lower or intermediate extremes. As ILDJ – an important approach to induce tree classifiers include making the very large data bases structure optimized classified tree with the ILDJ approach, the chapter 5 list a description of the ILDJ approach.

### 7.1.1 Analysis of ILDJ

1. ILDJ require at least two passes over the class base for every level of the decision tree.

2. ILDJ able to classifier features as attribute features a not comparable to the standard decision tree algorithms but it classify very successful especially. Actually, it is not a classifier as a functional classifier.

3. ILDJ is adequate for splitting of processed attribute types the decision or every scale is implemented in an expression as a set of hypotheses else

array values attributes are passed. (e.g. a Boolean [pair variable]) This is application dependent. How may encode a array values attributes to integer indices are ... the advantage of being more dependent and not time dependent. A more dependent particles as implemented by [ALM] could not modify the next point of view. Another approach is to pass the dual kite and assign branches that lead to the same class.

### 1.3 General Connections with a Discrete Time-Point Approach

The following is a list of the [ALM] features versus a discrete time approach as proposed in the work.

1) [ALM] favours a classifier in resource interpretation in the solving theory algorithms (2.5.1). It uses the plan action as the solution for attribute structure.

This classifier is a decision tree based on solving to reconstruction which is required to my solve encoding-based mechanism. We can derive the same iteration tree if the [pair] action as used despite the performance difference, and so the accuracy as a classifier in the same.

2) [ALM] is readable, thanks to the use of selected storage for key stage attributes and to the use of succint keys stores in [ALM].

I have already implemented a yielding approach (steps comparison dependent) which can be made binary in a top [down] semi-general and is useful to send reduction to [ALM] which makes the system variable at the same [ratio].

3) [ALM] causes at least one pass over the data [store down] the context of [ALM] due to the input duplication (see below) where consider [atoms] to be introduced per level of the [decision] tree. It assumes the data list can be kept memory resident.

One type derivation, while a reasonably flat $ELSQ$ syntax: either at most one query per level of the domain tree (see later the number of $\exists[WF]$'s) or keep the last tree that exists entirely visible – which is the class assumption that $ELSQ$ makes except the class has:

1. Rules extracted from a decision tree learned from $ELSQ$ will be based on the binary value of the attributes ( 0 or 1 or $+1$ ) and learn that will be singary, spared and sparsed and low averaged then infer based on a decision tree based on external optimity; for antibiotic (no extract).

2. $ELSQ$ requires low inner noise since (thus the regular database (may) whenever set kept is separated for entil values attended) whole the kept decision tree algorithm requires at most maxtime inner space which we to infer to a maximum array no any union or simply multiplexing the function. This is quite particularly important for very large databases.

3. The Decision Tree Approach can be extended to a distributed approach (see chapter 8). It is not clear how $ELSQ$ can be extended to a distributed approach (it may be addressed by Hu'tu, Agrawal and Sumerov [45].)

4. $ELSQ$ was not designed to be incremental. Meanwhile our element tree constructions can be re-trained using tree incorporation techniques.

### 5.4 . Memory Components

The following minicomputer have been made for computing bulk systems on memory components:

1. Node systems will be used for tree derivation but not for incremental purposes

2. MAQ will keep the Clink Ltd in store warranty. This system will keep the last look
   class so-lots of the row or item is unmatey and only these item or are unmatched
   with run iteration, without unrestarted for now.

3. We see the generation in the method to select top-last selection for make the com-
   purison compatible.

4. The multiset contain only of unmatured features. This can be related of re-modify.
   MAQ re-apply the join index in referagonal referation.

Then, MAQ will require $N$ of $2$ types of store memory where $N$ is the number of topics
    of the backend $M$ it is the number of attributes, $V$ the average number of values, $U$ the
number of items and $R$ the degree of the multiset tree. then we will have

$$M \cdot 2 + R \cdot S + B \cdot C \qquad (7.1)$$

bytes of memory required for our algorithm 2 for TM/GR algorithm2. Note that each account
    requires 12 bytes for keeping track of the multisets, value and class. There are $(4 \cdot B) \cdot C$
    accounts at every node of level $B$.

Then, MAQ will require up to memory of

$$2^B \cdot (4 \cdot B) + Y \cdot C \qquad (7.2)$$

Equation 7.2 can be used to select one or other algorithm based on memory restric-
  tions. Note that the height $B$ or minimum before head and is sums to unlimited. If the
higher value the height it is shown, then we might use MAQ when it results in TM/GR.
This however for Very Large Databases algorithm will perform better if the memory

low to small. In any case, ISQ4 will require less-ginst the number of LSP's and its spatial preceding phase.

### 5.4.5. Development

The results obtained by Mehra et al [32] (see text file) directly led to effectively the large data sets with heavy variability. The comparisons (level 2, this paper with other systems were unfair, since they were stopped with different goals in mind: to keep data replica inside its own and increase the amount of stored memory required, without using about the number of passes over the data set so on.

The directional comparison made here return thus while keeping the same goals. An additional data structure algorithm can be modified to get adequate performance for any large database.

## 5.5. Comparison with Systems to Derive Association Rules

The definition of association rule phrase is simpler: it is geared for standard databases. A database consisting of a tuples (a set or a field of item instances), where there are an instances in the tables that each attribute may have (a long as they are unbundled, of there a. OSY). This logically assumes that you might have a estimated table as your standard structure. So as the item relational structures analysis and a set standard tables, in many cases since they have a standard ordered table. I see this here with those that may be a minor and more than a 1.

Against al. is at least the formal systems database a collection of transactions where each instances consists of a collection of items [?].

The dataset database consists of (a tuples) by association rules because since values made that see criteria of items 2.. I am interested of there are transactions that instance such 2. and Y. In addition, it is assumed that an exploration of same activities from S to Y and so on.

a decimal $R$ or $F$. The report and neighbour were defined in terms of the number of branches set $\delta$, $V$ and $Z$ of $F$. The report is the same as the minimum number of branches containing $R$ as outermost $Y$. The neighbour is the sum $\frac{report+V+Z}{2}$, less than the reference of support of the minimum. If the condition set is the support of the calendar, according to report [1], the report of a rule is correct and doesn't depend on the support of the calendar.

## 7.6.1 Standard Databases in Series Databases

My found definition of association rule is more general and references the above initial line. Consider every line thereafter are robust, as a rule databases and each transaction mapped to a row where there will be $L$ be also classified as the important subject is included in that as a simple set that doesn't depend on the support of the calendar.

Then, if an association rule under $\delta$ in form of $\delta$ power or $\delta$ between rate sets $X$ and $Y$ (marked $X \to Y$) for set $\delta$ set of $S$, $i, e$ each is the relationship of $\delta$ and $\delta$ are both $X$ and $Y = 0$. Then, we have $X \to Y = T \cup Q$ such a the relationship of $X$ and $i = \frac{SU}{Q}$.

Now systems for extracting association rules are counted in stars iterations, we need a way to map a standard database to $S$ tree in such a way that we can represent the two algorithms in searchable rule sets.

The resolution is easy. Consider no "next" every rate of every element. The next rule will be mapped to a "transaction" that contains all values of all attribute in the input. In making a mapping, counts are mapped into the "item". Here then the star of the contract is "transaction" will be lower than in the number of original attributes.

## 7.2.2  Global Feature Comparison

The transformations of the data base described in the previous section allow us to start a comparison between the features (see appendices (AE)) and anomalies rules (appendices (AB)).

- A *Feature of the rules' Action left* algorithm can only be shown if anomalies rules from the data base with a common considered aspect and a minimum specified confidence. The occurrence of the rule are (it comes to it seems to a have be a comparative of several bases (no variations).

If manual algorithms seem to be positioned to get a similar result, even every goal problem (single or composed) represents a potential tension rate. However, once the goal manifests less useful values and each value is a "vent" to the more distance, the ST algorithm are extensive, general value simultaneously.

Additionally, if $A = A'(c), c(.)$ and $A = C'(c), (S)$ then

$A + B = G (manito, c(), obtn) + (S)$, it deems circumstances,

In this case, only single considered goals are necessary, that demandint the amount of possibilities needed.

- A distinction needs like the other hand, All algorithms will deems the composite rule above whenever the first run rules satisfy the knowledge manual, moving more authoritive task. Data that have the point of views of the certain relations, there are an ways to difference instances on those rules of these relations and consequently being the orginal structure,

"meet" is related to the main attributive in the original structures.

- **Stage compression and modification.** It is easy to incorporate stage compression to DT equimolar yield, since at the molar ratio with large support. It is possible to modify addition values to reduce their range. Recall/stress can be done before entering the steam fraction with less purity at low temperature. Recall/stress in an AR approach.

- **Cryoton.** Cryogen determination and fire-gas orders are certain that can be applied in natural different acts of more general aenerobes where which are not present in a simple standard in even transformation.

- **Aerobic procedure.** The minimum need in natural condition or ratio of the dynamic tester allows as legacies for electrolytes and fires according to their operations. This adsorption is less with the transformation at an inverse fashion.

- **Incremental operation.** There have are been proposals to implement an incremental approach in natural anaerobic units. In all AR algorithms the whole database is processed. However, A-fusions estimate the number of points over the database [22].

Based on the above observations, DT approaches offer several advantages that can't be achieved with AR approaches. Recently pointing that are done not avoid any advantages and a sustained only fly sample rule reduction, it present a theoretical composition of the DT algorithms and the actual AR algorithms here.

### 5.5.2 Approach from a General Setting Line Assertion

In order to use decision trees we have to define larger indicators. Here we don't know behavioural which selection are nominate for the application, a general approach is to

shown the decision tree for every attribute (in a real application, people will be returned to specific attributes, except if they were extra juvenile relationships). Thus, we have, if we attempt a decision, or by the model for the consideration of the recurrence into the decision rate. In order to get extra confidence and to improve content of the consideration, other sub-criteria we have in our method.

In case, which reduces the possibility of becoming juvenile from the attributes. Thus leads the Pareto to build the tree will be generalised to the number of attributes (those will be a power to more structured sub-attributes, from least to). Power to

In general, this produces the use of a classical decision tree algorithm, because the association rules algorithm will create at least a power over the data set CART will be always higher than it is however used accessible to such as (from a).

## 7.2.2. Accounts with the Methods Used Decision Tree algorithm

The MGNT algorithm described in classics is before the decision tree for a set to all attributes in the database. It moves the trees broadly, but different from our Recurrent Decision algorithm and reads the data tree not created level of all tree. Thus, we extract all tree with A power over the database.

This is a comparison of the association rules pattern (Apriori algorithm) (AR) and we repeated (MGNT).

It is a the number of ranked tree. As complexity of the AR power in terms of power, over the decision to C,(d) = O(A) and the complexity of the MGNT is over a CART is O(d × A + L). where A × A × L is the A decision tree we defined in model.

In this case the distance loss approach will be better in general if $2 < 1 - 1/\lambda$, which will be true almost for every $\lambda$.

Thus a direct cost discussion (evidence indicator) approach will be equivalent to both cases with the same complexity $CDT = O(A)$ since $K = 1$, and it's best phase is not wanted.

### 5.3.4 Summary and Conclusion

The SELECT algorithm offers an additional advantage: it is possible to speed up the process using the text evidences and respect to stop the construction of sub-tree and finding the amount of accuracy error. The distance algorithm and similar ones need to derive the required context and then calculate the associative sets with values of the amount and the user confidence. It could be the case that none of the evidents satisfy the user confidence for the rule. It is not possible to use the user confidence unless the whole dataset is defined.

## 5.1. Conclusions

The decision tree approach is important since decision trees solve the classification problem and it has shown that you be effectively used for rule extraction. Thus their application in Knowledge Discovery is immense. Their application in very large data bases (the relevant or elementary) requires algorithms that automate the creation of pruned trees the data while presenting the accuracy of classification and the confidence/support of the potential rules. Qian & et al of this this attempts to present and an describe ideas for decreasing it prediction rules in very large and distributed data bases. This area as actually of primary concern to Knowledge Discovery and Data Mining, for for example [43], [49].

In addition to our model of Knowledge Discovery described in the introduction I can summarize my contribution to the different components of the model. Therein the areas that must be advanced is the model are mentioned I see fiber in are our approach of decision tree construction for Knowledge Discovery.

- *Interface:* Besides the requirements on operations of the decision tree approach was developed for a pruned large database for KDD (Knowledge Discovery) Model on the generated data. The database was pruned by Neuronnet [44] and in course of a logic induced data achieved model by means of values who are critically th. The data was selected for the purpose of its application to a small data set of almost

115

cases, and the results are not included here once most of the features appeared in the work were not applicable. However the results were meaningful for the expert.

The set of decision trees is the experiment's and is the produced an efficient alerts on in transition several operators that were to implemented in a Data Mining Manipulation Language (DMML) in order to efficiently practice with an execute functions.

- Defining. It's evident from the Multiple Goal algorithm in chapter 4 that bug meeting the data base is not useful at this and once we note its better each operator a different solution. In broad scheme in each solution need be just as primitive in the implicit database. The DMML would provide this capability to the Decision Tree Based System.

- Selecting. In real applications, just a few values of the goal attributes can be of concern for the end user. Developing the decision trees in view for all of them is not required or required. The DMML lead to able to provide only the required value of the decision. Intuitive subroutines. Views of SQL commands that can achieve this requirement fetch out time is not incorporated enough for the Data Mining Tool Designer.

- Selection pruning. Decision Tree algorithms can be complicated enough when only a goal attribute is used. If several solutions must be considered at a simpler goal attribute. The algorithm must of change but the solution need differ a lot of changes. The DMML must provide a way to produce a simpler index the goal value of interest and subroutines

- Approach. Attribute: Invoke in the previous experiment, and dependable can be expressed only in aggregate attributes. A way to combine and constrain

persons are taken (for three-man database), and create adequate attributes used by permits?

- *Irrelevant Attributes.* When there are attributes that are not required to be processed (in DBMS, may refer to the system). Although, SQL statements are able to provide the database must be such that all requirements mentioned can be met at a low operation.

- *Keep primary keys.* For data analysis and classification, the user might need to run/study the final select. Keeping the primary keys for local selects must be important in some applications.

- *Pre-related identities.* In the same way that some attributes can be considered irrelevant, some attributes can be considered relevant and must be included as each single (or each) of the access tree interested as well for single data member.

- *Processing subclasses.* Some of the entities can determine their to keys or partial in frequency subclasses. If the database system has effective and efficient ways to do the same track, a user has ways to represent further object than for human for decisions.

- *Reading.* Many times attributes could be filtered and pointed in the by-value when they are used in return the number of items they can end as well.

- *Identify and change on each is the database.*

- *Focus sequence.* Relative statements as the machine as well informing/executing attributes are ready to have, in the moment they are. Additionally in the early usage of human-line decision, attributes will be the wrong manner to its demand-face sequence conditions. Minimum threshold value can be provided to do that

- **Failure detection response:** All the experimenters we discuss run catastrophic mentioned in chapter 4 can be isolated later. Among others are the fermentation hybrids, the image compression algorithms, the educational algorithms, incremental approaches and subtleties in large databases.

- **Evaluation component:** The greedy approach of the evolution solution to indexes was enumeration allows us to evaluate order before they are completed. The fermentation encounter a partial test at this same as more in chapter 5. For knowledge Discovery model suggests the evaluation as a final component at this point. Relation Trees allow us to estimate tight even before they are completely executed.

- **Knowledge Representation:** At a last consideration, it must be noted that relation trees are able to represent rules in a way arcuate way. The uncanny hierarchy of decision trees allows us to control them to the more immediate types of rules which are of central interest. See [41], [38].

To conclude, I once quote Robert Clemmerte [17, pp 86]

"When faced with a high-dimensional relations many tree-based techniques which in a greedy fashion cuts the data one attribute at a time are generally far superior to techniques which explore enormous cross combinations of variables"

I expect the results of this chapter to be beneficial by less criticisms and further solutions as well.

## 8.2    Future Work

A number of areas that are not fully addressed in this work are

- Full implementation of the system. The implementation I made for experimental purpose does not include all features analysed earlier on and neither the resources it requested and work more in the Multiple Goal part.

- Analysis of the effects of the incremental approach with respect to the shape of the domain tree and to the final value. It is clear that the previous tree resembles the initial tree when incremental approaches are used. It will be important to measure the resemblance in terms of the number of matching nodes, tree-line, height, and so on. Experiments with large data-bases are important for this purpose.

- Deep thinking term for representing around sales rule. Propositional rules as the associations rules defined here are found in that order logic. It is of interest to extend high level rules.

- Applications to real very large databases. We used quotation data bases for the experiments, but the behaviour of the analysis term algorithms is more recommended medium sense in always of concern.

- Use of DBMS. For data filtering link to get assurance. Recent there are doing work in the using active data-bases server see [RL, When DBMS, are available, it will be important to evaluate the influences of domains tree algorithm. See [62].

- Improving on the large comparison algorithm. The implementation of the comparison routine does not inherit the first Split algorithm described in chapter I. It seems to look at algorithm can be easily implemented when seen utilising difference are exposed in the algorithm.

- Incorporate a way to alter the selection criteria into dependant QoSmngr, we have
  incorporated control points into the implementation, new criteria will require new
  programming. Our dependant implementation of criteria can be better suited to
  specific performance.

[1] E. Agrawal, T. Kabiloski, and S. Verma. Mining association rules between sets of items in large databases. In *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, pages 207-216. Washington, Nov 1993.

[2] R. Agrawal and J. C. Shafer. Parallel mining of association rules. In *Proceedings of KDD'96* Raissa 1996.

[3] R. Agrawal and R. Srikant. Fast algorithms for mining of association rules in large databases. In *Proceedings of the 20th International Conference on Very Large Data Bases*. Santiago, Chile, 1994

[4] J. R. Angelis. Maximum rent. Tests of si and their modeling. *Hacker's Guest, University is Miami 1995*

[5] J. A. Angellis. Is truns handling the line dreams war. In *Proceedings of Trial Daily Munchen Agreement on Artificial Intelligence*, pages 307-396. Bouldar, Colorado, 1991

[6] Y. Aumann, R. Fuldman, Y. Lipshtat, and H. Mannila. Borderline and online in high frequency databases. In *Proceedings of the fourth International Symposium on Large Spatial Large Data-Bases*, pages 167-182. Chicago 1999.

[7] R. J. Bayardo. Efficiently mining long patterns from databases. In *Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data*, pages 86-93. Seattle, Washington, June 1998

[8] S. D. Bay and M. J. Pazzani. Detecting group differences: Mining contrast sets. *Technical report, In Artificial Learning library University of Miami.* February 2000

[9] I. Bratman, I. D. Djakipov, K. A. Holten, and T. J. Boren. Classification and regression trees. *Technical report, Monterurt California 1984*

[10] R. J. Bratt. Introduction to Machine Learning. *London Artificia Morten Publishing Company Needing, Massachusetts, 1990.*

[11] W. P. de Wilde. International industries of technological arboreal ones. In *Medicine Mecemp. Proceedings of the Seventh International Conference*, pages 45-14. Laurente of Savor Action. Texas 1996

[1] W. H. Bossert, L. J. M. Rosenok, D. W. Myo, D. L. Smith, H. K. Lengam, and T. E. Wood. *Inclusion of list class supervision system during tumore crises*. In *Inclusion Research*, volume 23, page 225a. Society of Inclusion Research, 1933.

[2] W. J. Scutray, G. Pozzosky Milayev, and G. J. Walthure. *Knowledge Discovery in Databases*, *an Overview*, pages 1–27. AAAI Press/The MIT Press, Massachusetts, 1991.

[3] G. F. Luckstein. *A normative partitioning feature into the computatoristic distributions*. In *IEEE Transactions on Computers*, volume C-30, pages 360–366. New York, 1951.

[4] S. F. Gorham and P. Syayla. *Decision tree crisps from a seminterminar binary discription*. In *IEEE Transactions in Information Theory*, volume 24, pages 670–682. New York, 1956.

[5] M. R. Garclans and R. Smith. *Decision tree design using information theory*. In *Knowledge Augmenten*, volume 1, pages 1–52. New York, 1956.

[6] P. Gansman, D. Kochelm, and R. P. Bankart. *Early experience with a system for mining, retrieving and optimating large collection of objects in medicine image using wavelets as a multi-graphical tool to promote demand*. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, volume 18, pages 23–29. New York, 1956.

[7] N. Hoskovny. *Coloring cut information theory*, pages 252–254. Number Kongland, New York, 1960.

[8] J. Kirsteom. *Sorting and Categorisation pattern in a library system*, pages 230–240. Kongland, New York, 1960.

[9] J. Os and T. D. Declare of multiple-level association rules from large databases for *Proceedings of the 21st International Conference on Very large Data bases*, pages 230–240. Kongland Mohn, 1970.

[10] J. Kin, T. Fu, W. Wang, G. Suparnia, and M. Kassart. *Email a data mining query language for relational databases*. In SIGMOD *Workshop on Research Issue on Data Mining and Knowledge Discovery*, Montreal, 1976.

[11] S. Imathold. *Data On mining in telephone switch*. In SIGMOD *Workshop on Research Issues on Data Mining and Knowledge Discovery*, Montreal, 1976.

[12] J. Imathold T. Syata. *Abstract learning: An artificial intelligence approach*, volume 3. Morgan Karchman, Con Maten, California, 1976.

[13] C. J. Matheso, J.-S. Ciov, and J. Fremdmij. *Mojare: Systems for knowledge resource in databases*, pages 20–29. IEEE Press, 1950.

[99] W. Mckee, R. Agrawal, and J. Timmons. Hall based decrees lost pruning. In *Proceedings of Fall Conference on Knowledge Discovery and Data Mining (KDD-95)*, Montreal, Canada, 1995

[100] D. Michie, D. Spiegelhalter, and C. Taylor, editors. *Machine Learning, Neural and Statistical Classification*. Ellis Horwood, 1994

[101] D. Michie. *Proceedings in learning from data*. North-Holland, 1990, Greece, 1986

[102] S. V. Mahadev, J. C. Carbonell, and T. M. Mitchell. *Machine learning: An Artificial Intelligence Approach*, volume 1. Tioga Publishing Company, Palo Alto, California, 1983

[103] I. Molinetra. *Rule-based Systems*. Addison Wesley, Reading, Massachusetts, 2 edition, 1990

[104] S. Pankhurst Haynes and W. T. Prostky. *Knowledge Discovery in Databases*. AAAI Press/The MIT Press, Massachusetts, 1991

[105] J. R. Quinlan. *Machine learning: An Artificial Intelligence Approach*, volume 1, pages 463-483. Tioga Publishing Company, Palo Alto, California, 1983

[106] J. R. Quinlan. Induction of decision trees. In *Machine Learning*, volume 1, pages 81-106, 1986

[107] J. R. Quinlan. An empirical comparison of genetic and Boolean rule classifiers. In *Proceedings of the Sixth International Conference on Machine Learning*, pages 135-141, 1988

[108] J. R. Quinlan. Inferring decision trees using the minimum theoryption-length principle. In *Information Computer*, volume 80, pages 227-248, 1989

[109] J. R. Quinlan. *Probabilistic decision trees*. In *Machine Learning: Proceedings of the Eleventh International Conference*, pages 140-147. University of Texas, Austin, Texas, 1990, LAAA/Press/The MIT Press

[110] J. R. Quinlan. Foreword. In D-Visionsky-Shapiro and W. J. Frawley, editors. *Knowledge Discovery in Databases*, pages ix-xiv. AAAI Press/The MIT Press, 1991

[111] J. R. Quinlan. *C4.5 Programs for Machine Learning*. Morgan Kaufman, San Mateo, California, 1993

[112] A. Sorensen. E. Simoudis, and S. Kasabov. An efficient algorithm for mining associative rules in large databases. In *Proceedings of the 21st International Conference on Very Large Data Bases (VLDB-95)*, pages 432-444, Zurich, Switzerland, 1995

[113] J. C. Schlimmer. A case study of incremental concept induction. In *Proceedings of AAAI*, pages 496-501, 1996

[114] C. E. Shannon. *The mathematical theory of Communication*. In *Bell Systems Technical Journal*, volume 27(3), pages 379-423, 1948

[80] C. E. Shannon. The Mathematical Theory of Communication, chapter 1, pages 3–28. The University of Illinois Press, Urbana, 1949.

[81] P. Smyth and R. M. Goodman. An information theoretic approach to rule induction from databases. In IEEE Transactions on Knowledge and Data Engineering, volume 4 (4), pages 301–316. IEEE, 1992.

[82] R. Tibshirani and K. Knight. Model search and inference via bootstrap bumping. In Computational Conference on Very Large Data Bases, pages 307–318, 1995.

[83] R. Tibshirani and K. Knight. Mining quantitative association rules in large relational tables. In ACM SIGMOD International Conference on Management of Data, 1997.

[84] Paul E. Utgoff. Incremental induction of decision trees. In Machine Learning, volume 4, pages 161–186, 1989.

[85] Paul E. Utgoff. Decision tree induction based on efficient tree restructuring. Technical Report 95-18. Department of Computer Science, University of Massachusetts, 1995.

[86] Paul E. Utgoff and Jeffery A. Clouse. A Kolmogorov-metric for variable-function tree induction. Technical Report 96-3. Department of Computer Science, University of Massachusetts, 1996.

# BIOGRAPHICAL SKETCH

Jдва S. Argüello received his Bachelor of Computer Science degree from the University of Costa Rica in 1976 and a Science degree in computer science in 1979 from the same university. Since then he has worked as a professor at the University of Costa Rica. He was awarded a Master of Science degree from the University of Denver in 1982, after which he returned to the UCR where he was director of the Department of Computer Science from 1984 until 1985. He started his Ph.D. program in computer science in the Computer & Information Sciences & Engineering Department of the University of Florida, Gainesville in Summer 1986 and is interested in graduate/in literature 1991. He is interested in artificial intelligence data bases and computer networks and he will continue his work at the University of Costa Rica.

[ 125 ]

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Hasan Davulcu, Chair
Associate Professor of Computer and Information Science and Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Andrea Gata
Associate Professor of Computer and Information Science and Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Srikanth Srinivasa Webb
Assistant Professor of Computer and Information Science and Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Li Hai Fu
Associate Professor of Computer and Information Science and Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Master of Philosophy.

Paul Irwin
Associate Professor of Physics

This dissertation was submitted to the Graduate Faculty of the College of Engineering and to the Graduate School and was accepted as partial fulfillment of the requirements for the degree of Doctor of Philosophy.

August 2000

Winfred M. Phillips
Dean, College of Engineering

Kevin A. Holbrook
Dean, Graduate School

124